



# Language-Based Engineering:

A Comprehensive Approach to Software Analysis and Hardening

**Hyunsoo Shin**

Program Coverage

Buffer Overrun:  
Address Sanitizer

Fuzzer

Symbolic  
Execution

Data Race  
Detector

Delta Debugging



LLM-based Fuzzer

First Print  
08-25-2025

Copyright © 2025 Hyunsoo Shin. All rights reserved.

PUBLISHED BY HYUNSOO SHIN

(“ <https://github.com/flavio20002/typst-orange-template> “, “TEMPLATE-WEBSITE”)

This book contains or is derived from software licensed under the Apache License, Version 2.0. You may obtain a copy of the License at: <http://www.apache.org/licenses/LICENSE-2.0>

Unless required by applicable law or agreed to in writing, software distributed under the License is distributed on an “AS IS” BASIS, WITHOUT WARRANTIES OR CONDITIONS OF ANY KIND, either express or implied. See the License for the specific language governing permissions and limitations under the License.

This book uses code from the “orange-template” project, stated above, licensed under the Apache License, Version 2.0. Copyright 2025 “orange-template” Authors.

No part of this publication may be reproduced, distributed, or transmitted in any form or by any means without the prior written permission of the author, except for brief quotations used in reviews or academic references.

*First printing, August 2025*

Contents

I	Introduction to Language-based engineering	
1	Welcome aboard . . . . .	13
1.1	As a Programmer . . . . .	13
1.2	About the book . . . . .	13
1.2.1	Contents . . . . .	14
1.2.2	Source and Target Language . . . . .	14
1.2.3	LLVM Version . . . . .	14
1.2.4	Prerequisite . . . . .	14
1.2.5	Environment . . . . .	14
1.2.6	Compile and test . . . . .	15
1.2.7	Rationale . . . . .	15
2	Prerequisite: Instrumentation, Runtime, and LLVM IR . . . . .	17
2.1	Compiler Pass Exploting . . . . .	17
2.2	Runtime Library . . . . .	18
2.3	Benefit of IR . . . . .	19
2.4	LLVM IR Tutorial . . . . .	22
II	Coverage	
3	Instrumentation and Runtime . . . . .	27
3.1	Introduction to Coverage . . . . .	27
3.2	Hardhat Coverage Example . . . . .	27
3.3	Goal . . . . .	29
3.4	How To Instrument . . . . .	30
3.4.1	Instrument Strategy . . . . .	31
III	Buffer Overrun: Address Sanitizer	

<b>4</b>	<b>Introduction to Out-of-bound Access Bug . . . . .</b>	<b>49</b>
<b>4.1</b>	<b>Buffer Overrun . . . . .</b>	<b>49</b>
4.1.1	Goal . . . . .	51
4.1.2	Reference Paper . . . . .	52
<b>4.2</b>	<b>Idea . . . . .</b>	<b>52</b>
4.2.1	Red Zone . . . . .	53
4.2.2	Allocator Overloading . . . . .	53
4.2.3	Shadow Memory Translation . . . . .	54
4.2.4	Implementation . . . . .	56

## IV Fuzzing

<b>5</b>	<b>Fuzzer . . . . .</b>	<b>81</b>
<b>5.1</b>	<b>Introduction to Fuzzer . . . . .</b>	<b>81</b>
<b>5.2</b>	<b>A Big Picture . . . . .</b>	<b>83</b>
<b>5.3</b>	<b>Design and Architecture . . . . .</b>	<b>84</b>
5.3.1	Branch(Edge) Coverage . . . . .	84
5.3.2	Bucketing . . . . .	85
5.3.3	Fork Server . . . . .	86
5.3.4	Harness . . . . .	86
5.3.5	Oracle . . . . .	87
5.3.6	Fuzzing Campaign . . . . .	87
<b>5.4</b>	<b>Implementation . . . . .</b>	<b>87</b>
5.4.1	Instrumentation . . . . .	88
5.4.2	Runtime . . . . .	92
5.4.3	Fuzzer . . . . .	99
5.4.4	Summary . . . . .	124

## V Symbolic Execution

<b>6</b>	<b>Introduction to Symbolic Execution . . . . .</b>	<b>129</b>
<b>6.1</b>	<b>Problem . . . . .</b>	<b>129</b>
<b>6.2</b>	<b>Working Example . . . . .</b>	<b>130</b>
<b>6.3</b>	<b>Design . . . . .</b>	<b>131</b>
6.3.1	State . . . . .	131
6.3.2	Static Analysis vs. Dynamic Management . . . . .	133
6.3.3	Interface (Harness) . . . . .	133
<b>6.4</b>	<b>Implementation . . . . .</b>	<b>134</b>
6.4.1	Instrumentation . . . . .	134
6.4.2	Runtime . . . . .	146

## VI Delta Debugging

<b>7</b>	<b>Introduction to Delta Debugging . . . . .</b>	<b>159</b>
<b>7.1</b>	<b>Introduction . . . . .</b>	<b>159</b>
<b>7.2</b>	<b>Simple Idea (Binary Search) . . . . .</b>	<b>159</b>
<b>7.3</b>	<b>Local Minmum and Global Minimum . . . . .</b>	<b>160</b>

7.4	<b>1-minimal Algorithm</b> . . . . .	161
7.5	<b>Implementation</b> . . . . .	161
7.6	<b>Testing</b> . . . . .	164
7.7	<b>Fuzzer Integration</b> . . . . .	166
7.7.1	Homework . . . . .	169

## VII LLM-based Synthesis

<b>8</b>	<b>LLM-based Synthesis</b> . . . . .	<b>173</b>
8.1	<b>Hello, AI</b> . . . . .	173
8.2	<b>Introduction to LLM-based Fuzzer</b> . . . . .	174
8.2.1	Autoprompting . . . . .	174
8.2.2	Fuzzing Loop . . . . .	175
8.3	<b>Introduction to Program Synthesis</b> . . . . .	176
8.4	<b>LLM-based Synthesis</b> . . . . .	177
8.4.1	Homework . . . . .	184
8.5	<b>Retrospect</b> . . . . .	184

## VIII Data Race Detector

<b>9</b>	<b>Introduction to Data Race</b> . . . . .	<b>187</b>
9.1	<b>Introduction</b> . . . . .	187
9.2	<b>Theory</b> . . . . .	191
9.2.1	Happens-Before . . . . .	191
9.2.2	LockSet . . . . .	193
9.2.3	Implementation . . . . .	196
9.2.4	Homework . . . . .	224

## IX Appendix

<b>10</b>	<b>Conclusion</b> . . . . .	<b>229</b>
10.1	<b>Further Study</b> . . . . .	229
	<b>Bibliography</b> . . . . .	<b>233</b>



## List of Figures

2.1	Concept of Module and Runtime . . . . .	18
3.1	input = [1,2,3,4,5] . . . . .	28
3.2	input = [1,3,5,7,9] . . . . .	29
3.3	input = [2,4,6,8,10] . . . . .	29
3.4	input = [1,2,3,4,5] . . . . .	29
3.5	input = [1,3,5,7,9] . . . . .	29
3.6	input = [2,4,6,8,10] . . . . .	30
4.1	Shadow Memory Layout: Adds redzones to both the left and right of the original memory layout. . . . .	53
4.2	Cases: When OOB access is detected vs. not detected . . . . .	53
4.3	Boundary poisoning . . . . .	55
5.1	A naive fuzzer . . . . .	82
5.2	Advanced fuzzer . . . . .	82
5.3	A workflow of fuzzing campaign . . . . .	83
5.4	Implementation Components . . . . .	88
5.5	Implementation Components . . . . .	117
5.6	Fuzzer Execution . . . . .	124
6.1	Symbolic State . . . . .	133
6.2	Leaf State Nodes . . . . .	145
7.1	1-Minimal Algorithm . . . . .	161
8.1	Overview of Fuzz4All . . . . .	174
8.2	AutoPrompting . . . . .	175
8.3	Fuzzing Loop . . . . .	176
9.1	A and B send message to C, C sends to D. (ts = timestamp) . . . . .	191
9.2	Lockset Algorithm . . . . .	193
9.3	Example . . . . .	194
9.4	Example . . . . .	195





Buffer Overrun: Address  
Sanitizer

Symbolic Execution

Fuzzer

Delta Debugging

Program Coverage

LLM-based Fuzzer

Data Race Detector

List of Tables



## Part I

# Introduction to Language- based engineering

<b>1</b>	<b>Welcome aboard . . . . .</b>	<b>13</b>
1.1	As a Programmer . . . . .	13
1.2	About the book . . . . .	13
	<b>Prerequisite: Instrumentation,</b>	
<b>2</b>	<b>Runtime, and LLVM IR . . . . .</b>	<b>17</b>
2.1	Compiler Pass Exploring . . . . .	17
2.2	Runtime Library . . . . .	18
2.3	Benefit of IR . . . . .	19
2.4	LLVM IR Tutorial . . . . .	22



# printf(“Hello, World”);

## 1. Welcome aboard

### 1.1 As a Programmer

As usual, I was coding and Googling to find the traditional use cases of a topic I had just encountered for the first time. I clicked on one of the top search results—a blog—and on its front page, I saw a phrase written in Korean: “개발자로 살길 정말 잘했다” (“I’m really glad I chose to be a programmer”).

Yes, I had forgotten the joy of programming for some time. I realized I’d been too focused on building programs quickly and forcing myself to study academic topics. I set high standards for myself in an effort to become a great programmer, but it ended up making me accumulate knowledge under stress and pressure. No one was pushing me to do that — I’m not even sure why I was. Maybe it’s just in my DNA.

Writing this book turned out to be a stroke of luck. It brought back the fragrance of joy that programming once gave me. Honestly, I have to admit — while writing this book, there were times I was just focused on finishing it as quickly as possible. Looking back, I must have been a little crazy.

I truly hope this book brings you the same sense of relaxation and enjoyment I felt while writing it. Enjoy!

### 1.2 About the book

There are a lot of research area in computer science field such as Artificial Intelligence, Network, HCI (Human-Computer-Interaction), Security, PL (Programming Language), Computer Architecture, Operating System, Computer Graphics, Computer Vision, Robotics, Database, Distributed System, Cryptography, etc.

Personally, I don’t prefer research areas that require additional equipment. On the contrary, I’m drawn to fields where I can do meaningful work with just a laptop. That’s why programming languages (PL) really suit my personality — I never feel the need for extra resources to study it. However, I sometimes struggle with PL because it’s deeply rooted in algorithms, mathematics, and logic. I’ve seen many brilliant minds in this field, and I often find motivation in their work.

Secondly, I’ve been interested in the field of security. Security itself is a broad area, but my focus has been on language-based security [1]; — that is, studying security through

the lens of programming languages (PL). I’m also particularly interested in topics like optimization, abstraction, and correctness.

This book explores a series of common program bugs that many of us have likely encountered while coding. At its core, the goal is to identify bugs — from a security perspective — **using PL technologies**.

### 1.2.1 Contents

This book delves into seven major topics related to program analysis and software testing:

#### Topics

1. Coverage
  - a. Reports line, branch, and function-level coverage during program execution.
2. Buffer Overrun: Address Sanitizer
  - a. Detects out-of-bounds memory accesses on the heap and stack.
3. Fuzzing
  - a. Develops a coverage-guided fuzzer inspired by AFL (American Fuzzy Lop).
4. Delta Debugging
  - a. Minimizes crashing inputs generated during fuzzing by isolating the minimal failure-inducing input.
5. Symbolic Execution
  - a. Solves branch conditions to enhance code coverage during fuzzing.
6. LLM-based Synthesis
  - a. Explore LLM-based fuzzer and implement LLM-based synthesis to generate test code snippets
7. Data Race Detection
  - a. Identifies data races in multi-threaded programs.

### 1.2.2 Source and Target Language

The implementations in this book are written in Rust. Theoretically, the target language is the C language family. During the writing of this book, only C was thoroughly tested, with limited experimentation in C++. However, the approach can be extended to any language that supports an LLVM IR backend.

Broad language support is achieved through the frontend, which is responsible for recognizing the intermediate representation (IR) and providing the corresponding handlers within the instrumentation module.

### 1.2.3 LLVM Version

The book uses LLVM version 17.0.6, which was released on November 28, 2023.

### 1.2.4 Prerequisite

Readers are expected to be familiar with Rust programming. This book does not cover Rust syntax, coding conventions, or other language fundamentals.

### 1.2.5 Environment

We provide a Nix environment that sets up LLVM 17 in your shell. Simply type `nix-shell` in the top-level directory where `default.nix` is located. This environment also includes all the dependencies required for the book.

### 1.2.6 Compile and test

The build recipe is defined in the justfile located in the top-level directory.

- `just build`: Compile the entire project
- `just test`: Run unit test

### 1.2.7 Rationale

The code examples provided in this book may contain bugs, typographical errors, or inaccurately explained concepts. If you identify any issues, contributions via pull requests are welcome and appreciated.





# Three Principles



## 2. Prerequisite: Instrumentation, Runtime, and LLVM IR

### 2.1 Compiler Pass Exploding

Human writes a program using high-level syntax like:

```
1 int sum = a + b;
```

C

Compiler converts this high-level expression into lower expression like:

```
1 %a = load i32, ptr %1
2 %b = load i32, ptr %2
3 @sum = add i32 @a, @b
```

LLVM-IR

Last lowering is transformation into bytecode, which is machine independent such as x86 and AMD.

```
1 100010101010010...
```

Bytecode

From a programmer's perspective, the compilation process—from intermediate representation (IR) to bytecode—is typically abstracted away. Developers focus on writing high-level source code, which serves as input to the compiler. For instance, when attempting to measure code coverage, manually inserting function calls to track line execution is both tedious and error-prone. Moreover, this approach does not scale well, as it would require repeating the same effort for each program individually.

This is where instrumentation comes into play, introduced here for the first time in this book. An instrumentation tool automatically inserts function calls (e.g., for line hits) into the source or IR, relieving the programmer of manual effort. This process is typically performed within the compiler's middle layer. For example, Clang [2]; supports provides an instrumentation framework based on LLVM Passes [3].

In this book, we implement instrumentation entirely in Rust as a modular system. The module takes IR as input, performs transformations, and outputs instrumented IR. Each instrumentation module can target a different concern, and multiple instrumentation layers can be applied sequentially—for example, instrumenting first with Module A, then B, and finally C.

A sample of instrumented IR is shown below:

```

1 %a = load i32, ptr %1
2 %b = load i32, ptr %2
3 @sum = add i32 @a, @b
4 call line_hit(...) // instrumented

```

LLVM-IR

## 2.2 Runtime Library

Compilers generate executable code, which can be run directly on a CPU. However, there are inherent limitations to what can be achieved solely through statically generated code.

For example, memory allocation in C is performed using `malloc()`, a POSIX-standard interface [4]. The actual memory management is handled by the runtime library (e.g., `glibc`), which internally invokes system calls to the operating system. In modern programming languages, garbage collection automates memory management, relieving the programmer from manual allocation and deallocation. Like `malloc()`, garbage collectors are part of the runtime system, not the compiler's output.

Go (Golang) is a prominent example of a language that supports reflection [5]. At runtime, Go programs can inspect type information, examine generic values such as `interface{}` or `any`, and iterate over struct fields. These capabilities are enabled by the runtime library, not the compiler alone.

In summary, runtime libraries play a crucial role in extending the dynamic behavior of programming languages beyond what the compiler can provide.

In this book, we implement a custom runtime library as a shared library written in Rust. It interacts with the instrumented code during execution, enabling dynamic analysis and runtime features.

Having now introduced the two core components of this book—**instrumentation** and **runtime**—we're ready to dive into implementation.

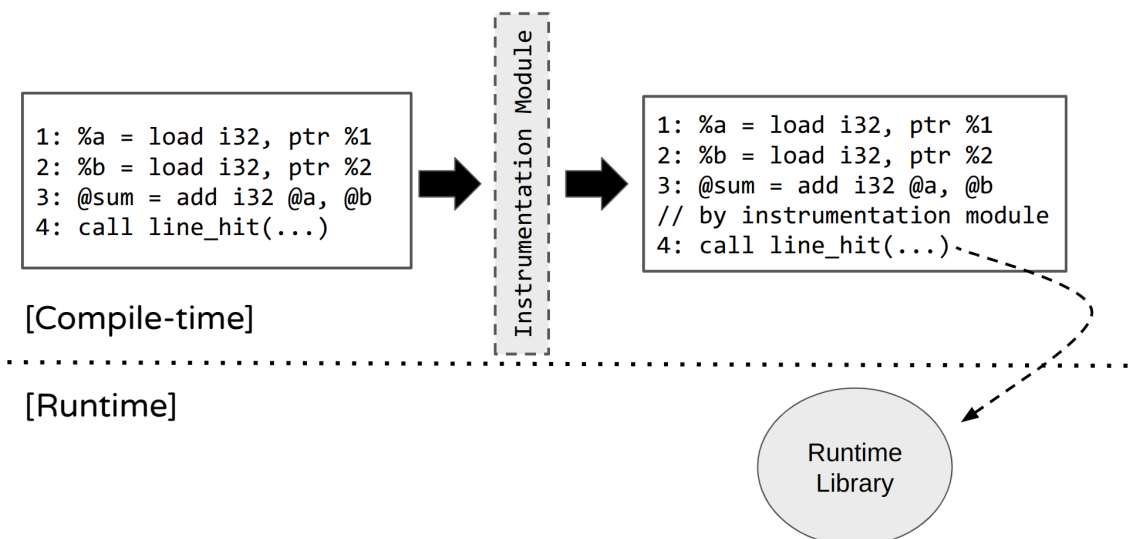


Figure 2.1: Concept of Module and Runtime

## 2.3 Benefit of IR

This chapter introduces the principles of IR, with a focus on LLVM IR. Before diving into IR itself, let's first ask: *Why is IR well-suited for program instrumentation?* Beyond instrumentation, IR is also a common target for tasks such as program analysis, correctness checking, and formal verification.

Let's take a look example C code.

```
1  #include <stdio.h>
2
3  #define ARR_SIZE 5
4
5  int get_odd_even_diff(int* arr, int size) {
6      int even_cnt = 0;
7      int odd_cnt = 0;
8      for (int i=0; i<size; i++) {
9          if (arr[i] % 2 == 0) {
10             even_cnt++;
11         } else {
12             odd_cnt++;
13         }
14     }
15     return odd_cnt - even_cnt;
16 }
17
18 int main() {
19     int arr[ARR_SIZE] = {1,2,3,4,5};
20     printf("%d\n", get_odd_even_diff(arr, ARR_SIZE));
21 }
```

The function `get_odd_even_diff()` takes a pointer to an integer array along with its size. It counts the number of odd and even elements in the array, then returns the difference between the count of odd and even numbers. For example, given the array {1, 2, 3, 4, 5}, the function returns 1 because there are 3 odd numbers and 2 even numbers, and  $3 - 2 = 1$ .

Below is the corresponding LLVM IR code for the `get_odd_even_diff()` function.

```
1  define i32 @get_odd_even_diff(ptr noundef %0, i32 noundef
   %1) #0 !dbg !16 {
2      %3 = alloca ptr, align 8
3      %4 = alloca i32, align 4
4      %5 = alloca i32, align 4
5      %6 = alloca i32, align 4
6      %7 = alloca i32, align 4
7      store ptr %0, ptr %3, align 8
```

```

8    call void @llvm.dbg.declare(metadata ptr %3, metadata !22,
   metadata !DIExpression()), !dbg !23
9    store i32 %1, ptr %4, align 4
10   call void @llvm.dbg.declare(metadata ptr %4, metadata !24,
   metadata !DIExpression()), !dbg !25
11   call void @llvm.dbg.declare(metadata ptr %5, metadata !26,
   metadata !DIExpression()), !dbg !27
12   store i32 0, ptr %5, align 4, !dbg !27
13   call void @llvm.dbg.declare(metadata ptr %6, metadata !28,
   metadata !DIExpression()), !dbg !29
14   store i32 0, ptr %6, align 4, !dbg !29
15   call void @llvm.dbg.declare(metadata ptr %7, metadata !30,
   metadata !DIExpression()), !dbg !32
16   store i32 0, ptr %7, align 4, !dbg !32
17   br label %8, !dbg !33
18
19 8:                                     ; preds = %27, %2
20   %9 = load i32, ptr %7, align 4, !dbg !34
21   %10 = load i32, ptr %4, align 4, !dbg !36
22   %11 = icmp slt i32 %9, %10, !dbg !37
23   br i1 %11, label %12, label %30, !dbg !38
24
25 12:                                    ; preds = %8
26   %13 = load ptr, ptr %3, align 8, !dbg !39
27   %14 = load i32, ptr %7, align 4, !dbg !42
28   %15 = sext i32 %14 to i64, !dbg !39
29   %16 = getelementptr i32, ptr %13, i64 %15, !dbg !39
30   %17 = load i32, ptr %16, align 4, !dbg !39
31   %18 = srem i32 %17, 2, !dbg !43
32   %19 = icmp eq i32 %18, 0, !dbg !44
33   br i1 %19, label %20, label %23, !dbg !45
34
35 20:                                    ; preds = %12
36   %21 = load i32, ptr %5, align 4, !dbg !46
37   %22 = add i32 %21, 1, !dbg !46
38   store i32 %22, ptr %5, align 4, !dbg !46
39   br label %26, !dbg !48
40
41 23:                                    ; preds = %12
42   %24 = load i32, ptr %6, align 4, !dbg !49
43   %25 = add i32 %24, 1, !dbg !49
44   store i32 %25, ptr %6, align 4, !dbg !49

```

```

45    br label %26
46
47 26:                                ; preds = %23, %20
48    br label %27, !dbg !51
49
50 27:                                ; preds = %26
51    %28 = load i32, ptr %7, align 4, !dbg !52
52    %29 = add i32 %28, 1, !dbg !52
53    store i32 %29, ptr %7, align 4, !dbg !52
54    br label %8, !dbg !53, !llvm.loop !54
55
56 30:                                ; preds = %8
57    %31 = load i32, ptr %6, align 4, !dbg !57
58    %32 = load i32, ptr %5, align 4, !dbg !58
59    %33 = sub i32 %31, %32, !dbg !59
60    ret i32 %33, !dbg !60
61 }

```

The most striking aspect of IR is its fine-grained nature compared to high-level source code. For example, the condition `if (arr[i] % 2 == 0)` appears as a single line in high-level code, but in IR, it's broken down into multiple instructions:

- Access of array index operation
- A modulo operation
- An equality comparison
- A conditional branch

When working with high-level code, we would need to manually decompose and parse each operator for such fine-grained analysis.

Another key feature of IR is that variables are never reassigned. In LLVM IR, an instruction like `%1 = add i32 %a, %b` introduces a new variable `%1`. If another computation follows, it generates a new variable, such as `%2`, rather than reusing `%1`. This property is known as **Static Single Assignment (SSA)** form [6]. SSA greatly simplifies program analysis by making it easier to track the origin and use of each value. Let's look at an example.

```

1  int a = 1; // scope 1
2  if (...) { // scope 2
3      int a = 2;
4      printf("%d\n", a); // the variable in line 3 of scope 2 is
                          referenced here.
5  }

```

Without SSA, tracking variable scopes and updates becomes our responsibility. However, with SSA, this burden is lifted—we can rely on the SSA form to handle variable scoping automatically. As a result, program analysis becomes significantly simpler and more reliable.

```

1  int a_1 = 1;
2  if (...) {
3      int a_2 = 2;
4      printf("%d\n", a_1);
5  }

```

Lastly, IR explicitly represents control flow. In the `get_odd_even_diff()` function, for example, if the for loop condition evaluates to false, execution jumps to line 11. Similarly, if the if condition on line 5 is false, control transfers to line 8. LLVM IR uses the concept of basic blocks, which are sequences of instructions with no internal jumps or branches—control only enters at the beginning and exits at the end. All high-level control structures (if, for, while, etc.) are translated into a series of basic blocks connected by branch instructions. Thanks to this structure, we don't need to manually determine the next program counter (PC); the IR makes control flow explicit and analyzable.

In summary, three key benefits highlight the importance of focusing our business logic on instrumentation, analysis, and related tasks.

- ✓ Decomposed instructions
- ✓ Control-flow
- ✓ SSA

## 2.4 LLVM IR Tutorial

First, to generate IR file, we should add additional compiler flags `-S` and `-emit-llvm`. To generate debug information, attach `-g` as well. For example, The above IR example was generated by `clang -g -O0 -S -emit-llvm get_odd_even_diff.c -o get_odd_even_diff.ll`.

Function is defined as follows:

```

1  define i32 @get_odd_even_diff(ptr noundef %0, i32 noundef
   %1) #0 !dbg !16 {
2      ...
3      %4 = alloca i32, align 4
4      ...
5      br label %8, !dbg !33
6
7  8:                                     ; preds = %27, %2
8      ...
9      br i1 %11, label %12, label %30, !dbg !38
10
11 20:                                     ; preds = %12
12      ...
13      %21 = load i32, ptr %5, align 4, !dbg !46
14      store i32 %22, ptr %5, align 4, !dbg !46
15      ...

```

- The `define` keyword is used to declare a function. The return type `i32` indicates that the function returns a 32-bit integer. The argument (`ptr noundef %0, i32 noundef %1`) means the function takes a pointer type and a 32-bit integer parameter named `%0` and `%1`, respectively. The `noundef` qualifier indicates that the argument must not be an undefined value.
- Line 3 allocates memory for a 32-bit integer (`i32`). At this point, it's not specified whether the integer is signed or unsigned. Later instructions, such as `srem i32 %a, %b` or `urem i32 %a, %b`, determine how it is treated.
- Line 5 transfers control to the basic block labeled `#8`. The next instruction to be executed is the first instruction of the basic block `#8`.
- Line 9 performs a conditional branch based on the value of `%11`. If the condition evaluates to true, control is transferred to basic block `#12`; otherwise, it goes to basic block `#30`.
- Line 13 loads a `i32` type of value from pointer `%5` and stores it into `%21`.
- Line 14 stores the value of `%22` into the memory location pointed to by `%5`. This stored value can be retrieved later using a load instruction.

As mentioned earlier, a basic block is a straight-line code sequence with no branches except at the end. In cases where the instruction sequence becomes long, it will break up into multiple basic blocks, even if there is no explicit control transfer at each step.

#### Remark 2.1

#### Summary

All seven prepared topics operate on the IR. The instrumentation modules insert IR code at specific target locations. Using the IR binding library, we can create and manipulate global variables, functions, basic blocks, and instructions — including the ability to remove instructions.





## Part II

# Coverage

	<b>Instrumentation and Runtime. . .</b>	
<b>3</b>	<b>27</b>	
3.1	Introduction to Coverage . . . . .	27
3.2	Hardhat Coverage Example . . . . .	27
3.3	Goal . . . . .	29
3.4	How To Instrument . . . . .	30



# “How much of your code is tested?”

32% → 67% → 100%

## 3. Instrumentation and Runtime

### 3.1 Introduction to Coverage

What is the coverage [7]? Coverage is a measurement of which lines of code have been executed.

Have you ever used a coverage tool? Coverage is most commonly used in unit testing. Many developers aim to increase coverage to gain confidence in how thoroughly their code is being tested. For example, 100% coverage means that all lines of code are executed during testing, while 70% coverage indicates that 30% of the code has not been exercised by the tests.

Code coverage is one of the key topics covered in Software Engineering courses. Several global companies provide coverage tools that can be integrated with your CI pipelines or third-party tools like IDEs.

The goal of this chapter is to develop a **coverage instrumentation module** and a corresponding **runtime library** to measure function, branch, and line coverage. The implementation will be done in Rust, targeting programs written in C.

To guide our development, we'll model our tool after Hardhat's coverage tool [8]. Hardhat is a smart contract development framework that simplifies testing and deployment. Achieving high coverage in smart contracts is critical, as vulnerabilities can lead to real financial loss. From my experience, Hardhat Coverage is one of the most intuitive coverage tools I've used, which is why I chose it as a reference model.

### 3.2 Hardhat Coverage Example

First, let's take a look at a Hardhat Coverage example. The semantics of the example code are the same as those of the `get_odd_even_diff` function shown in this example.

```
1 // SPDX-License-Identifier: UNLICENSED Solidity
2 pragma solidity ^0.8.28;
3
4 contract TestContract {
5     function getOddEvenDiff(int256[] memory arr) public pure returns
      (int256) {
```

```

6      int256 evenCnt =0;
7      int256 oddCnt =0;
8      for (uint i=0; i<arr.length; i++) {
9          if (arr[i] % 2 == 0) {
10             evenCnt++;
11         } else {
12             oddCnt++;
13         }
14     }
15     return oddCnt - evenCnt;
16 }
17 }

```

Let's write a test for the contract.

```

1  const { expect } = require("chai");
2
3  describe("Test", function () {
4      it("Test", async function () {
5          const factory = await ethers.getContractFactory("TestContract");
6          const contract = await factory.deploy();
7
8          expect(await contract.getOddEvenDiff([1,2,3,4,5])).to.equal(1);
9      })
10 });

```

Once you run `npx hardhat coverage`, the test code will be executed, and a summary will show which lines were hit and which were not, like this:

File	% Stmts	% Branch	% Funcs	% Lines	Uncovered Lines
contracts/ C.sol	100 100	100 100	100 100	100 100	
All files	100	100	100	100	

Figure 3.1: input = [1,2,3,4,5]

The input [1,2,3,4,5] covers all lines of code. Now, let's try changing the input to [1,3,5,7,9] and [2,4,6,8,10].

```

1  describe("Test", function () {
2      it("Test", async function () {
3          ...

```

```

4     expect(await contract.getOddEvenDiff([1,3,5,7,9])).to.equal(5);
5     ...
6   })
7 });

```

The input [1,3,5,7,9] did not cover the line 10.

File	% Stmts	% Branch	% Funcs	% Lines	Uncovered Lines
contracts/ C.sol	100 100	50 50	100 100	85.71 85.71	10
All files	100	50	100	85.71	

Figure 3.2: input = [1,3,5,7,9]

The input [2,4,6,8,10] did not cover the line 12.

File	% Stmts	% Branch	% Funcs	% Lines	Uncovered Lines
contracts/ C.sol	100 100	50 50	100 100	85.71 85.71	12
All files	100	50	100	85.71	

Figure 3.3: input = [2,4,6,8,10]

### 3.3 Goal

Finally implemented our coverage will report for the function `get_odd_even_diff` is as follows:

File	% Funcs	Uncovered Funcs	% Branch	Uncovered Branches	% Lines	Uncovered lines
get_odd_even_diff.c	100.00		100.00		100.00	

Figure 3.4: input = [1,2,3,4,5]

File	% Funcs	Uncovered Funcs	% Branch	Uncovered Branches	% Lines	Uncovered lines
get_odd_even_diff.c	100.00		75.00	10 (12:F)	86.67	10, 11

Figure 3.5: input = [1,3,5,7,9]

File	% Funcs	Uncovered Funcs	% Branch	Uncovered Branches	% Lines	Uncovered lines
get_odd_even_diff.c	100.00		75.00	12 (10:T)	93.33	12

Figure 3.6: input = [2,4,6,8,10]

You can compile and run as follows:

```

> clang -g -O0 -S -emit-llvm get_odd_even_diff.c -o
1 get_odd_even_diff.ll // compile origin program and produce IR
2 > ./instrument -i get_odd_even_diff.ll -o iget_odd_even_diff.ll -m
  coverage // instrument the compiled IR
3 > clang iget_odd_even_diff.ll -L. -lcoverage_runtime // compile the
  instrumented IR and produce binary
4 > LD_LIBRARY_PATH=./ ./a.out // run the binary

```

There are seven columns in our coverage report. The table below describes each reported item.

Reported item	Description
File	File name
% Funcs	Degree of function coverage
Uncovered Funcs	Uncovered function lines
% Branch	Degree of branch coverage
Uncovered Branch	Uncovered branch lines
% Lines	Degree of line coverage
Uncovered Lines	Uncovered lines

There are two conditional branches—one for the for loop and one for the if-else statement inside it. If the loop condition evaluates to true, the body executes; otherwise, control jumps outside the loop. This results in a total of four possible branches. For the inputs [1, 3, 5, 7, 9] and [2, 4, 6, 8, 10], the for loop branches are fully covered, but either the if or else branch inside the loop is never taken. As a result, only 3 out of 4 branches are executed, leading to 75% branch coverage.

### Remark 3.1

### Summary

We introduced one of the most promising coverage tools, Hardhat Coverage. Our goal is to implement a similar UI and functionality. Additionally, we will provide detailed reports showing uncovered lines for each function, branch, and line.

## 3.4 How To Instrument

We will implement the initial instrumentation module to collect three types of coverage data: function, branch, and line coverage.

### 3.4.1 Instrument Strategy

First, we need to distinguish between functions, branches, and lines in the LLVM IR.

Most LLVM IR binding libraries provide utility functions to iterate over functions defined in a given IR module. Inkwell—the LLVM IR binding library [9] used by the book—also offers such functionality. This allows us to retrieve the definition line of each function by querying its debug metadata. To access this metadata, the IR must be generated with the `-g` compilation flag, which includes debugging information such as line numbers.

Second, for identifying branch lines, we must inspect each instruction to determine whether it is a branch instruction (***br ...***) or a switch instruction (***switch ...***). By iterating over the instructions within each basic block and checking their opcodes, we can accurately identify these control flow instructions.

Lastly, all instructions within a basic block is directly associated with lines.

#### 3.4.1.1 Iterating IR file

First, the given IR file is parsed into a module structure managed by the binding library. Our analysis begins with this module. The module contains functions, each function contains basic blocks, and each basic block contains instructions—forming the hierarchy: “IR file → Module → Function → Basic Block → Instruction.”

The following code demonstrates our common pattern for iterating over an IR module.

```
1 let funcs: Vec<_> = module.get_functions().collect();
2 for func in funcs {
3     for basic_blk in func.get_basic_blocks() {
4         for instr in basic_blk.get_instructions() {
5             ...
6         }
7     }
8 }
```

#### 3.4.1.2 Collecting line numbers

We need to collect all line numbers associated with functions, branches, and individual lines. The following function retrieves the line number for a given LLVM value, which may represent either a function or an instruction.

```
1 // instrument/src/llvm_intrinsic.rs
2 pub fn get_instr_loc<'ctx, T: AnyValue<'ctx>>(instr: &'ctx T) ->
   (u32, u32) {
3     unsafe {
4         let line = LLVMGetDebugLocLine(instr.as_value_ref());
5         let col = LLVMGetDebugLocColumn(instr.as_value_ref());
6         (line, col)
7     }
8 }
```

Here is how we collect line numbers.

```

1 // instrument/src/coverage.rs Rust
2 let module_filename = cstr_to_str(module.get_source_file_name());
3 let mut brs_loc = vec![];
4 let mut lines_loc = BTreeSet::new();
5 let mut lines = BTreeSet::new();
6 for instr in basic_blk.get_instructions() {
7     // Instrument within only given IR.
8     // When we work with C++ IR, a lot of unknown files are
9     // collected together
10    // We're interested for only our target C/C++ files, not other
11    // language system files.
12    if let Some(instr_filename) = get_instr_filename(&instr) {
13        if instr_filename != module_filename {
14            continue;
15        }
16    }
17    // Get line number of instructions and collect into a set.
18    let (line, _) = get_instr_loc(&instr);
19    lines.insert(line);
20    lines_loc.insert(line);
21    // Get line number of branch instruction and collect into a set.
22    record_br(&mut brs_loc, &instr);
23    record_switch(&mut brs_loc, &instr);
24 }

```

The handling of branch instructions is defined as follows:

```

1 // instrument/src/llvm_intrinsic.rs Rust
2 pub fn record_br(brs_loc: &mut Vec<u32>, instr: &InstructionValue) {
3     // if instruction is conditional branch
4     if instr.is_conditional() {
5         // if is `invoke` instruction
6         if instr.get_opcode() == InstructionOpcode::Invoke {
7             // get labels of normal and exception
8             let (next_normal_label, exception_label) = (
9                 instr.get_operand(0).unwrap().right().unwrap(),
10                instr.get_operand(1).unwrap().right().unwrap(),
11            );
12            if let Some(mut first_instr) =
13                next_normal_label.get_first_instruction() {
14                let loc = get_br_loc(&mut first_instr);
15                brs_loc.push(loc);
16            }
17        }
18    }
19 }

```



```

16         if let Some(mut first_instr) =
           exception_label.get_first_instruction() {
17             let loc = get_br_loc(&mut first_instr);
18             brs_loc.push(loc);
19         }
20     }
21     // if is `br` instruction
22     if instr.get_opcode() == InstructionOpcode::Br {
23         // get labels (basic blocks) of true and false branch
24         let (tbr, fbr) = (
25             instr.get_operand(2).unwrap().right().unwrap(),
26             instr.get_operand(1).unwrap().right().unwrap(),
27         );
28         // collect line numbers
29         if let Some(mut first_instr) =
           tbr.get_first_instruction() {
30             let loc = get_br_loc(&mut first_instr);
31             brs_loc.push(loc);
32         }
33         if let Some(mut first_instr) =
           fbr.get_first_instruction() {
34             let loc = get_br_loc(&mut first_instr);
35             brs_loc.push(loc);
36         }
37     }
38 }
39 }

```

The invoke instruction typically appears in C++ code. Consider the example below:

```

1 try {
2     ... // normal label
3 } catch (...) {
4     ... // exception label
5 }

```

CPP

If a runtime error occurs—such as an out-of-bounds access—within the code inside the try statement, control is transferred directly to the exception label. Otherwise, the exception label is never executed.

### 3.4.1.3 Build instrumentation

Once we have gathered the corresponding line numbers for a given IR, the final step is to instrument runtime functions that flag which lines have been executed.

A naive implementation is to instrument every IR line individually, such as:

```

1 %1 = call ...
2 // call cov_hit_line(...)
3 %2 = load ...
4 // call cov_hit_line(...)
5 %3 = add ...
6 // call cov_hit_line(...)

```

LLVM-IR

The original IR consists of only three lines. In a naive instrumentation approach, a runtime function call is inserted for each IR line, resulting in a total of six lines—effectively doubling the code size.

Instead of instrumenting each line individually, we can optimize this by inserting a single runtime call per basic block. Since we already collect the line numbers for each basic block (as shown above), we declare a local array containing these line numbers.

As explained in [Benefit of IR](#), all instructions within a basic block are guaranteed to execute before any control transfer (e.g., a branch) occurs. Therefore, we gather the line numbers of all instructions within a basic block, store them in a local array, and pass this array to a single runtime function call.

Below is the transformed version of the [get\\_odd\\_even\\_diff](#) function. The key changes include:

- Declaring a local array of line numbers
- Making a runtime function call with this array

```

1  define i32 @get_odd_even_diff(ptr noundef %0, i32 noundef
   %1) #0 !dbg !9 {
2      %__cov_hit_lines_arr = alloca [5 x i32], align 4 // [instrumented]
   array declaration
3      store [5 x i32] [i32 0, i32 5, i32 6, i32 7, i32 8], ptr
   %__cov_hit_lines_arr, align 4 // [instrumented] intialize array
4      %__cov_hit_lines_arr_ptr = getelementptr [5 x i32], ptr
   %__cov_hit_lines_arr, i32 0, i32 0 // [instrumented] get pointer
   of the arr
5      %3 = call i64 @__cov_hit_batch(ptr @for.c, ptr
   %__cov_hit_lines_arr_ptr, i64 5) // [instrumented] call runtime
   function call with the array pointer
6      %4 = alloca ptr, align 8
7      %5 = alloca i32, align 4
8      %6 = alloca i32, align 4
9      %7 = alloca i32, align 4
10     %8 = alloca i32, align 4
11     store ptr %0, ptr %4, align 8
12     call void @llvm.dbg.declare(metadata ptr %4, metadata !15,
   metadata !DIExpression()), !dbg !16
13     store i32 %1, ptr %5, align 4
14     call void @llvm.dbg.declare(metadata ptr %5, metadata !17,
   metadata !DIExpression()), !dbg !18
15     call void @llvm.dbg.declare(metadata ptr %6, metadata !19,
   metadata !DIExpression()), !dbg !20

```

LLVM-IR

```

16  store i32 0, ptr %6, align 4, !dbg !20
17  call void @llvm.dbg.declare(metadata ptr %7, metadata !21,
18    metadata !DIExpression()), !dbg !22
18  store i32 0, ptr %7, align 4, !dbg !22
19  call void @llvm.dbg.declare(metadata ptr %8, metadata !23,
20    metadata !DIExpression()), !dbg !25
20  store i32 0, ptr %8, align 4, !dbg !25
21  br label %9, !dbg !26

```

Lines 1 to 3 are instrumented—they declare an array, initialize it, and invoke a runtime function. No additional instrumentation is inserted within this basic block. In the next basic block, the same process is repeated. Therefore, we insert three lines of instrumentation code per basic block, rather than per individual line of code.

The following code snippets demonstrate how declarations and runtime function calls are constructed. The term **build** is commonly used in LLVM IR bindings and reflects the standard approach to IR instrumentation.

```

1  // instrument/src/coverage.rs Rust
2  let lines_len = lines.len(); => usize
3  if lines_len > 0 {
4      let arr_ptr = build_i32_static_arr(
5          context,
6          builder,
7          &lines.into_iter().collect(),
8          COV_HIT_LINES_ARR,
9          COV_HIT_LINES_ARR_PTR,
10         )?;
11     build_cov_hit_batch(
12         context,
13         module,
14         builder,
15         &filename_str_ptr.unwrap(),
16         arr_ptr,
17         lines_len,
18     )?;
19 }

```

If there is a line within the current basic block, install the local array and the runtime function call using the `build_i32_static_arr()` and `build_cov_hit_batch()` build functions, respectively.

```

1  // instrument/src/inkwell_intrinsic.rs Rust
2  pub fn build_i32_static_arr<'ctx>(
3      context: &'ctx Context,
4      builder: &Builder<'ctx>,

```

```

5     vals: &Vec<u32>,
6     alloca_name: &str,
7     gep_name: &str,
8 ) -> Result<PointerValue<'ctx>> {
9     let arr_typ =
10    context.i32_type().array_type(vals.len().try_into()?); // Define
    array type
11    let int_vals: Vec<_> = vals
12    .iter()
13    .map(|v| context.i32_type().const_int(*v as u64, false))
14    .collect(); // Gather element values (line numbers)
15    // Build(instrument) array
16    let arr_val = context.i32_type().const_array(&int_vals);
17    let arr_alloca = builder.build_alloca(arr_typ, alloca_name)?;
18    builder.build_store(arr_alloca, arr_val)?;
19
20    // Build declaration of array pointer
21    // Get a pointer to the first element
22    let zero = context.i32_type().const_int(0, false);
23    let array_ptr = unsafe {
24    builder.build_gep(
25    arr_typ,
26    arr_alloca,
27    &[zero, zero], // Get pointer to the first element
28    gep_name,
29    )
30    }?;
31    Ok(array_ptr)
32 }

```

To build a function call, we should create a function first.

```

1 // instrument/src/inkwell_intrinsic.rs 🔒 Rust
2 fn get_cov_hit_batch_func<'ctx>(
3     context: &'ctx Context,
4     module: &Module<'ctx>,
5 ) -> FunctionValue<'ctx> {
6     // If the function had been built before, get the function value
7     // Otherwise, build it.
8     match get_func(module, COV_HIT_BATCH) {
9         Some(func) => func,
10        None => {

```

```

11         let record_func_cov_lines_typ =
12             context.i64_type().fn_type(
13                 context.ptr_type(AddressSpace::default()).into(),
14                 context.ptr_type(AddressSpace::default()).into(),
15                 context.i64_type().into(),
16                 ],
17                 false,
18             );
19         module.add_function(COV_HIT_BATCH,
20                             record_func_cov_lines_typ, None)
21     }
22 }
23
24 // instrument/src/inkwell_intrinsic.rs
25 pub fn build_cov_hit_batch<'ctx>(<
26     context: &'ctx Context,
27     module: &Module<'ctx>,
28     builder: &Builder<'ctx>,
29     filename_str_ptr: &GlobalValue,
30     arr_ptr: PointerValue,
31     arr_length: usize,
32 ) -> Result<()> {
33     let record_func_cov_lines = get_cov_hit_batch_func(context,
34                                                         module);
35     // Call coverage runtime function with three parameters, file
36     // name, array pointer, and size of array.
37     builder.build_call(
38         record_func_cov_lines,
39         &[
40             filename_str_ptr.as_pointer_value().into(),
41             arr_ptr.into(),
42             convert_to_int_val(context,
43                                 arr_length.try_into()?).into(),
44         ],
45         "",
46     )?;
47     Ok(())
48 }

```

The process of creating a runtime function signature (function declaration) is extensively used in other modules as well.

So far, we do not know how many lines of code are covered. To obtain this information, the runtime must know both how many lines of code are executable and which lines have been executed. Currently, we have only handled the latter.

To calculate the coverage percentage, the final step is to build a runtime initialization function. This function collects all executable lines related to functions, branches, and lines within the IR module, not the basic block level. Later, this information is used to determine how many lines have been covered by comparing it with all executable lines. This build function is executed at the final phase of the instrumentation process.

```

1 // instrument/src/coverage.rs
2 let init_last_instr = get_cov_init_last_instr(module);
3 builder.position_before(&init_last_instr);
4 build_src_mapping_call(
5     context,
6     module,
7     builder,
8     &filename_str_ptr.unwrap(),
9     &funcs_loc,
10    &brs_loc,
11    &lines_loc,
12 )?;
```

Record the total number of code lines for functions, branches, and lines. First, a local array is constructed; then, this array is passed to a runtime function that initializes and records all the code line information.

```

1 // instrument/src/inkwell_intrinsic.rs
2 pub fn build_src_mapping_call<'ctx>(
3     context: &'ctx Context,
4     module: &Module<'ctx>,
5     builder: &Builder<'ctx>,
6     filename_str_ptr: &GlobalValue,
7     funcs_loc: &BTreeSet<u32>,
8     brs_loc: &Vec<u32>,
9     lines_loc: &BTreeSet<u32>,
10 ) -> Result<()> {
11     let funcs_loc_vec: Vec<_> =
12         funcs_loc.clone().into_iter().collect();
13     let brs_loc_vec: Vec<_> = brs_loc.clone().into_iter().collect();
14     let lines_loc_vec: Vec<_> =
15         lines_loc.clone().into_iter().collect();
16
17     let func_lines_ptr = build_i32_static_arr(
```

```

16     context,
17     builder,
18     &funcs_loc_vec,
19     COV_SRC_MAPPING_FUNC_LINES,
20     COV_SRC_MAPPING_FUNC_LINES_PTR,
21 )?;
22     let brs_lines_ptr = build_i32_static_arr(
23     context,
24     builder,
25     &brs_loc_vec,
26     COV_SRC_MAPPING_BRS_LINES,
27     COV_SRC_MAPPING_BRS_LINES_PTR,
28 )?;
29     let lines_lines_ptr = build_i32_static_arr(
30     context,
31     builder,
32     &lines_loc_vec,
33     COV_SRC_MAPPING_LINES_LINES,
34     COV_SRC_MAPPING_LINES_LINES_PTR,
35 )?;
36
37     builder.build_call(
38         get_src_mapping_func(context, module),
39         &[
40             filename_str_ptr.as_pointer_value().into(),
41             func_lines_ptr.into(),
42             convert_to_int_val(context,
43             funcs_loc_vec.len().try_into()?).into(),
44             brs_lines_ptr.into(),
45             convert_to_int_val(context,
46             brs_loc_vec.len().try_into()?).into(),
47             lines_lines_ptr.into(),
48             convert_to_int_val(context,
49             lines_loc_vec.len().try_into()?).into(),
50         ],
51         "",
52     )?;
53     Ok(())
54 }

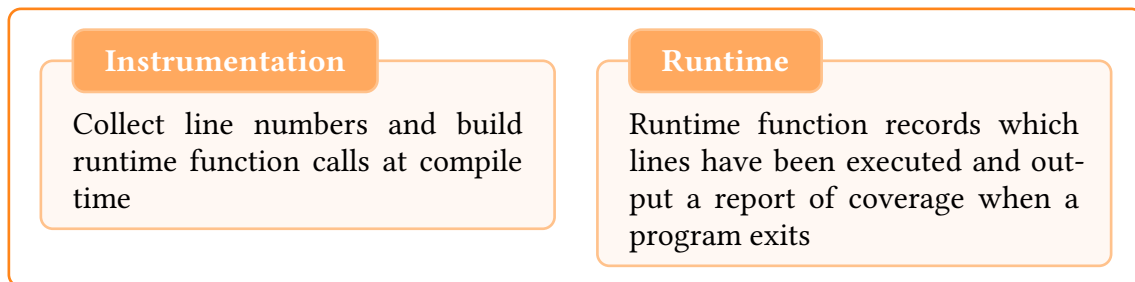
```

A key point is that the source mapping function is created within a **constructor** [10] function. In C/C++, a constructor function runs automatically when the process loads—before the `main()` function executes. If the given IR already defines constructor

functions, it's not an issue because multiple constructors can coexist and are prioritized based on a specified priority parameter. Our construction function works regardless of when it is executed within the set of construction functions.

#### 3.4.1.4 Coverage Runtime

We've instrumented a runtime function call for each basic block, and now it's time to define the runtime function itself. Before diving into the implementation, let's take a moment to recap the bigger picture of what we're trying to achieve.



There are three exposed runtime functions.

- `__cov_init()`: Called at the constructor function
  1. Register report function to be called at program exit via `libc::atexit()`
  2. Initialize internal coverage structure
- `__cov_mapping_src(...)`: Called at the constructor function
  - Record all code lines
- `__cov_hit_batch(...)`: Called per function, branch, and instructions
  - Record which lines are hit

```

1 // coverage_runtime/src/coverage_runtime.rs
2 #[no_mangle]
3 pub extern "C" fn __cov_mapping_src(
4     file_ptr: *const libc::c_char,
5     funcs_ptr: *const u32,
6     funcs_length: usize,
7     brs_ptr: *const u32,
8     brs_length: usize,
9     lines_ptr: *const u32,
10    lines_length: usize,
11 ) {
12     if file_ptr.is_null() || funcs_ptr.is_null() ||
13        brs_ptr.is_null() || lines_ptr.is_null() {
14         return;
15     }
16     let filename = cstr_to_string(file_ptr);
17     let (funcs, brs, lines) = unsafe {
18         (
19             std::slice::from_raw_parts(funcs_ptr, funcs_length),

```



```
20         std::slice::from_raw_parts(brs_ptr, brs_length),
21         std::slice::from_raw_parts(lines_ptr, lines_length),
22     )
23 };
24 let src_map = SourceMapping {
25     lines: lines.to_vec(),
26     brs: brs.to_vec(),
27     funcs: funcs.to_vec(),
28 };
29 COVERAGE_STATE
30     .write()
31     .unwrap()
32     .source_map
33     .lock()
34     .unwrap()
35     .insert(filename, src_map);
36 }
37
38 #[no_mangle]
39 pub extern "C" fn __cov_hit_batch(
40     file_ptr: *const libc::c_char,
41     lines_ptr: *const u32,
42     length: usize,
43 ) {
44     if lines_ptr.is_null() {
45         return;
46     }
47
48     // convert raw C pointer into rust slice
49     let lines = unsafe { std::slice::from_raw_parts(lines_ptr,
50         length) };
51     for line in lines {
52         __cov_record(file_ptr, *line);
53     }
54
55     /// Records a hit at the given source location
56     #[no_mangle]
57     pub extern "C" fn __cov_record(file_ptr: *const libc::c_char, line:
58         u32) -> usize {
59         // Early return if coverage is disabled
60         if COVERAGE_STATE
```

```

60     .read()
61     .unwrap()
62     .enabled
63     .load(Ordering::Relaxed)
64     == 0
65 {
66     return 0;
67 }
68
69 let file = cstr_to_string(file_ptr);
70 let loc = LineMapping { file, line };
71
72 // Get or create counter for this location
73 let lines_idx = {
74     let state = COVERAGE_STATE.write().unwrap();
75     let mut loc_map = state.location_map.lock().unwrap();
76     if let Some(&idx) = loc_map.get(&loc) {
77         idx
78     } else {
79         let idx = state.lines.len();
80         loc_map.insert(loc, idx);
81         idx
82     }
83 };
84 let mut state = COVERAGE_STATE.write().unwrap();
85 if lines_idx >= state.lines.len() {
86     // Insert a new index
87     state.lines.push(AtomicUsize::new(0));
88 }
89 // Increment the lines
90 state.lines[lines_idx].fetch_add(1, Ordering::Relaxed)
91 }

```

There are several ways to collect a report over the program's lifetime. In this book, we use the `atexit()` function [11] to register a report-generating function, which ensures that the coverage report is automatically generated when the program terminates.

The final function, `__cov_hit_batch()`, records hit information into the global state. This global state primarily manages the following two sets:

- `state.source_map`: Stores all code lines at runtime initialization
- `state.location_map`: Stores hit line per instructions

The `#[no_mangle]` annotation is necessary for our runtime functions to be callable from the instrumented program. Since we're writing the shared library in Rust, any functions we want to export must not be mangled.

If you've worked with C++, you may have encountered name mangling [12]. C++ supports function overloading, allowing multiple functions with the same name but different parameter types. To distinguish these during compilation, the compiler encodes type information into function names—a process known as name mangling. As a result, each overloaded version gets a unique symbol name.

In our case, the instrumented program must call the runtime function using its exact name. Therefore, name mangling must be disabled to ensure the function is exposed with the correct symbol name.

The report function computes coverage metrics at the function, branch, and line levels. Once this functionality is in place, we can improve the output presentation by formatting it into a table and adding colors for readability.

#### 3.4.1.5 Integration with Test Environment (gtest)

The most common use case for coverage information is its integration into unit testing environments. In this example, we'll use gtest, a highly popular C/C++ testing framework.

Navigate to the unit-test directory. Within this example, we'll simply demonstrate testing a basic math function:

```
1  #include <math.h>
2
3  int add(int a, int b) {
4      int v1 = a;
5      int v2 = b;
6      return v1 + v2;
7  }
8
9  int mul(int a, int b) {
10     int v1 = a;
11     int v2 = b;
12     return v1 * v2;
13 }
14
15 int unused(int a, int b) {
16     int v1 = a;
17     int v2 = b;
18     return v1 + v2;
19 }
```

math.c contains simple add, mul, and unused functions. The unused function will be detected as untouched in the coverage results.

This is the unit test file that tests math coverage:

```
1  #include <gtest/gtest.h>
2
3  extern "C" {
```

```

4     #include "math.h"
5 }
6
7 TEST(MathTest, Add) {
8     EXPECT_EQ(add(2, 3), 5);
9     EXPECT_EQ(mul(2, 3), 6);
10 }
11
12 int main(int argc, char **argv) {
13     ::testing::InitGoogleTest(&argc, argv);
14     return RUN_ALL_TESTS();
15 }

```

The test file exclusively invokes the add and mul functions. Consequently, the unused() function should be reported as uncovered.

Let's run

```

1  > just build
2  > just run
3  LD_LIBRARY_PATH=../../bin/debug ./test_math
4  [=====] Running 1 test from 1 test suite.
5  [-----] Global test environment set-up.
6  [-----] 1 test from MathTest
7  [ RUN      ] MathTest.Add
8  [          OK ] MathTest.Add (0 ms)
9  [-----] 1 test from MathTest (0 ms total)
10
11 [-----] Global test environment tear-down
12 [=====] 1 test from 1 test suite ran. (0 ms total)
13 [  PASSED  ] 1 test.
14 +-----+-----+-----+-----+
15 |   File | % Funcs | Uncovered Funcs | % Branch | Uncovered Branches |
16 | % Lines | Uncovered lines |
17 +-----+-----+-----+-----+
18 | math.c | 66.67 |          15 |      NaN |
19 | 69.23 | 15,16,17,18 |
20 +-----+-----+-----+-----+

```

As we can see, the function coverage is not 100%. All lines within the unused() function's body are uncovered, specifically reported as lines 15, 16, 17, and 18.

### 3.4.1.6 Summary

We've implemented the first instrumentation module and its corresponding runtime library. The instrumentation module operates at compile time, generating coverage-related information that calls the runtime library. The runtime library, in turn, manages and reports this coverage data during program execution. In essence, the instrumentation prepares the data we care about, and the runtime library processes it at runtime. This combination of compile-time and runtime techniques is a recurring pattern in the chapters that follow.

For the complete coverage code, refer to the `instrument` and `coverage_runtime` directories.

Coverage serves as a warm-up to help you become familiar with the instrumentation-runtime workflow. In the next chapter, we'll apply the same principle to build an out-of-bounds access detector known as AddressSanitizer (ASAN).

You can check out provided test cases in `coverage/tools/tests/inputs/coverage`.

### 3.4.1.7 Homework

**Exercise 3.1** We suggest the following topics for further study and exploration:

- **Topic #1: Reducing Runtime Function Calls with Dominance Graphs**
  - The current implementation places runtime function calls at the basic block level. However, this overhead can be further reduced by using a dominance graph [13]. By constructing a dominance relationship between basic blocks, we can identify opportunities to eliminate redundant instrumentation. For example, if the execution path is guaranteed to follow  $BB1 \rightarrow BB3 \rightarrow BB7$ , then placing the runtime function call in `BB1` alone is sufficient. This allows us to eliminate the calls in `BB3` and `BB7`, reducing runtime overhead.
- **Topic #2: Further information**
  - The runtime state variable `state.lines` tracks the number of executions per line. This information can be leveraged to identify hotspots or cold spots in the code. For example, lines with unusually high or low hit counts may indicate performance-critical code paths or dead code, respectively. Further analysis of this data can provide deeper insights into code behavior and optimization opportunities.



## Part III

# Buffer Overrun: Address Sanitizer

	<b>Introduction to Out-of-bound</b>	
<b>4</b>	<b>Access Bug . . . . .</b>	<b>49</b>
4.1	Buffer Overrun . . . . .	49
4.2	Idea . . . . .	52





# “Hardening the memory usage to improve safety”

RED ZONE

```
char* ptr =  
malloc(100);
```

RED ZONE

## 4. Introduction to Out-of-bound Access Bug

### 4.1 Buffer Overrun

Lists are among the most dominant data structures in modern programming languages. Many languages abstract raw arrays to provide additional functionality, such as dynamic sizing and built-in length tracking. It's hard to imagine programming without some form of list structure.

However, arrays (and lists) come with a well-known pitfall: out-of-bounds (OOB) access bugs. These occur when a program tries to access memory outside the valid range of an array or buffer.

Chances are, many readers have encountered OOB bugs at some point. Consider the following C code, which demonstrates an OOB access:

```
1  #include <stdio.h>
2
3  int main() {
4      int arr[3] = {1,2,3};
5      printf("%d\n", arr[2]); // 3
6      printf("%d\n", arr[3]); // ?
7  }
```

Line 6 accesses an index out of bounds. The `arr` occupies 12 bytes, allowing valid access only within the range of indices 0 to 2. Therefore, `arr[3]` accesses memory outside the defined bounds, resulting in undefined behavior. The returned value is unpredictable—some compilers may return zero, while others may produce garbage values. Regardless of the outcome, this is clearly a bug.

A more serious issue is when the program continues executing without reporting the bug. The invalid value may propagate through subsequent instructions, leading to undefined behavior (UB) [14]. UB makes debugging significantly harder.

Rather than allowing the program to continue silently, modern compilers and run-times often instrument checks to detect out-of-bounds (OOB) accesses and crash the program with a helpful call stack trace. As a result, modern programming languages

are generally better equipped to detect and report OOB bugs compared to traditional compilers like GCC or Clang.

Let's take a look the Golang code with the same sementic.

```
1 package main
2
3 import "fmt"
4
5 func main() {
6     arr := []int{1, 2, 3}
7     fmt.Println(arr[2]) // 3
8     fmt.Println(arr[3]) // crash
9 }
```

Go explicitly reports OOB errors and prevents the program from continuing execution. This behavior helps developers quickly identify and fix such bugs. When an OOB access occurs at runtime, Go outputs an error message as follows:

```
1 > go run oob.go
2 3
3 panic: runtime error: index out of range [3] with length 3
4
5 goroutine 1 [running]:
6 main.main()
7     .../oob.go:8 +0x87
8 exit status 2
```

Now, let's look at a Rust example.

```
1 fn main() {
2     let arr = vec![1,2,3];
3     println!("{}", arr[2]);
4     println!("{}", arr[3]); // crash
5 }
```

Compiling and running the code will produce the following output:

```
1 > rustc oob.rs && ./oob
2 3
3
4 thread 'main' panicked at oob.rs:4:23:
5 index out of bounds: the len is 3 but the index is 3
6 note: run with `RUST_BACKTRACE=1` environment variable to display a
  backtrace
```

### 4.1.1 Goal

As we observed the discrepancy between traditional compilers (*Gcc* and *Clang*) and modern compilers (*Go* and *Rust*), the latter have recognized common error patterns—such as OOB access—and made significant efforts to create a more robust programming environment through both instrumentation and runtime mechanisms. In this chapter, we will implement an OOB detector using a custom instrumentation module and a runtime library. **A naive C program will be compiled to LLVM IR, then instrumented using our custom ASAN instrumentation. If an OOB access occurs during execution, our ASAN runtime will detect and report it.**

For the **OOB code snipeet**, let's compile the C program and instrument, and attach runtime library.

```

1 > clang -g -O0 -S -emit-llvm asan-example.c -o iasan-
  example.ll // compile origin program and produce IR
2 > ./instrument -i iasan-example.ll -o iasan-example.ll -m asan //
  instrument the generated IR
3 > clang iasan-example.ll -L. -lasan_runtime // compile the
  instrumented IR and produce binary

```

Our OOB detector will report as follows:

```

1 > LD_LIBRARY_PATH=./ ./a.out // libasan_runtime.so should
  be located in the current directory. Otherwise you should
  specify the exact path where the library places.
2 3
3 [ASAN] invalid memory access detected at asan-example.c:
  0x7ffcfc3fbf20
4 0: asan_runtime::asan_runtime::report_asan_violated::{closure}}
5      at .../asan_runtime/src/asan_runtime.rs:46:18
6 1: std::thread::local::LocalKey<T>::try_with
7      at .../.rustup/toolchains/stable-x86_64-unknown-linux-
  gnu/lib/rustlib/src/rust/library/std/src/thread/
  local.rs:308:12
8 2: std::thread::local::LocalKey<T>::with
9      at .../.rustup/toolchains/stable-x86_64-unknown-linux-
  gnu/lib/rustlib/src/rust/library/std/src/thread/
  local.rs:272:9
10 3: asan_runtime::asan_runtime::report_asan_violated
11      at .../asan_runtime/src/asan_runtime.rs:44:5
12 4: __asan_mem_check
13      at .../coverage/asan_runtime/src/asan_runtime.rs:30:9
14 5: main
15      at ./asan-example.c:6:22
16 6: __libc_start_call_main
17 7: __libc_start_main_alias_2
18 8: _start

```

### Running the instrumented program with runtime library

The OOB detector reports a stack trace and identifies the location where the OOB access occurred. Lines 14 and 15 pinpoint the exact location of the issue. Specifically, in the main function, line 6 is identified as the point of the OOB access. Our detector is capable of handling both heap and stack buffer overflows.

#### 4.1.2 Reference Paper

In this chapter, we refer to a seminal paper that is considered a foundational work in this research area: “*AddressSanitizer: A Fast Address Sanity Checker*” (ASan) [15]. AddressSanitizer is a tool designed to detect out-of-bounds (OOB) memory errors. The paper primarily focuses on five categories of memory-related bugs. (Use After Free (UAF) [16], Buffer Overflow (BOF) [17], Uninitialized Read [16], Double Free [18], Null Dereference [19]).

ASan is currently built into the Clang compiler toolchain [20] and can be used in any C project.

Bug Category	Description
Use-After-Free (UAF)	Using memory after it's freed
Buffer Overflow	Accessing memory outside valid bound
Uninitializd Read	Reading memory without initialization
Double Free	Freeing memory more than once
Null Dereference	Dereferencing a null pointer

In this book, we implement only two categories of bugs: UAF and BOF.

#### 4.2 Idea

In the following code, the variable `arr` occupies 12 bytes of stack memory. To detect memory violations, shadow memory is introduced. Shadow memory tracks the allocation and deallocation status of each memory region, allowing the system to monitor when and where memory is allocated or freed.

```

1 int arr[3] = {1,2,3};
2 arr[2] // 3
3 arr[3] // 00B

```

C

### 4.2.1 Red Zone

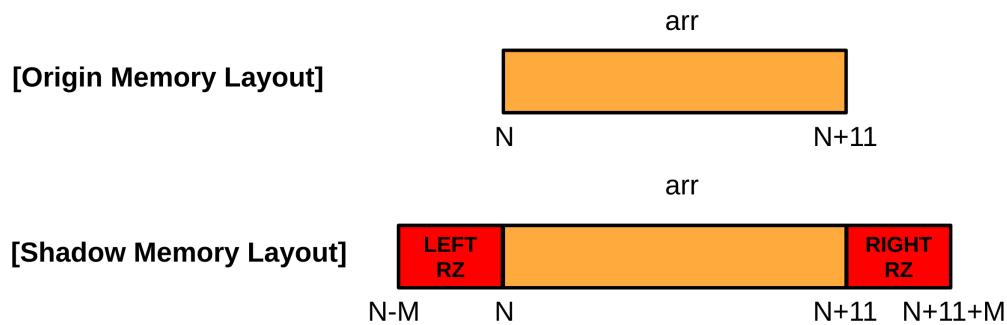


Figure 4.1: Shadow Memory Layout: Adds redzones to both the left and right of the original memory layout.

The figure illustrates two memory layouts: the original memory layout of the variable `arr` and its corresponding shadow memory layout. In the original layout, 12 bytes are allocated exactly as requested—nothing more. In contrast, the shadow memory layout includes additional space called “**redzones**” placed before and after the allocated region to detect OOB access. Assuming each redzone is 32 bytes, a total of 64 extra bytes are allocated—tagged as the left and right redzones, respectively.

Accessing index 3 is reliably detected because it falls within the right redzone, which is tracked by the shadow memory. However, if the access goes beyond the redzone boundaries, it may go undetected. This introduces a trade-off: increasing the redzone size improves the likelihood of catching out-of-bounds (OOB) accesses but also increases memory overhead. Conversely, reducing the redzone size optimizes memory usage but may lower the detection rate. Typically, a redzone size of 32 bytes is used as a balanced default. See the [two cases](#) illustration for examples.

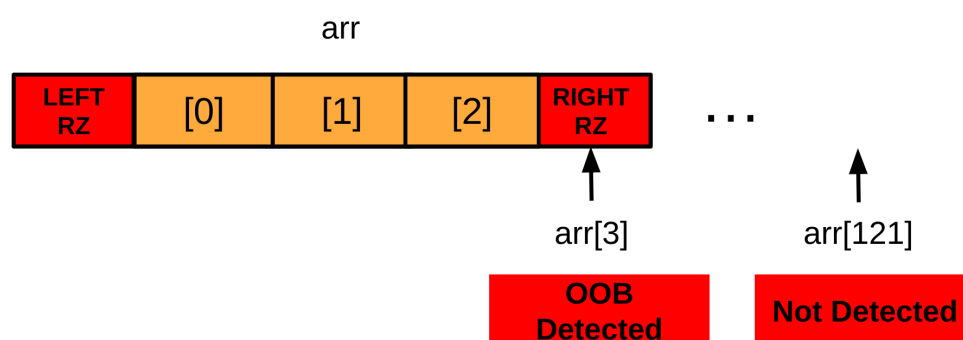


Figure 4.2: Cases: When OOB access is detected vs. not detected

### 4.2.2 Allocator Overloading

To automatically manage redzones, the ASan runtime must intercept memory allocation and deallocation requests. Here’s how it works:

- **Allocation** (e.g., 12-byte request):

- How many bytes are allocated?
  - A total of 76 bytes is allocated
    - 32 bytes for the left redzone
    - 12 bytes for the usable region
    - 32 bytes for the right redzone
    - Although the redzones are not used for program data, they must occupy memory to prevent the OS from reusing those regions.
- **Which address is returned?**
  - The allocator returns the address pointing to the start of the usable region (offset 32).
    - Address 0–31: Left redzone
    - Address 32–43: Usable memory
    - Address 44–75: Right redzone
    - The redzones are never returned to the program; they are not exposed through any legal pointer.
- **What happens in the shadow memory?**
  1. Shadow memory marks the left and right redzones with a special magic value to indicate that they are poisoned.
  2. During program execution, instrumented code checks the shadow memory to determine whether a memory access touches a poisoned (redzone) area.
- **Deallocation**
  1. The allocator frees the entire 76-byte block, including the redzones and the usable region.
  2. Shadow memory is updated to mark the entire region as freed.
- **Reallocation**
  1. Mark the entire existing region (including redzones) as freed in the shadow memory.
  2. Allocate a new block with redzones surrounding the requested size.
  3. Copy the contents of the original usable region into the new usable region.
  4. Deallocate the previous full block, including redzones.
  5. Return a pointer to the new usable region.

ASan allocates additional memory—called redzones—whenever a memory allocation request is made during program execution. It also maintains a separate region known as shadow memory, which is managed exclusively by the ASan runtime. This shadow memory is reserved when the ASan runtime library is first loaded.

In the referenced paper, shadow memory occupies one-eighth of the virtual address space. Allocating such a large memory region is not problematic due to how OSes manage memory. Specifically, virtual memory is reserved upfront, but physical memory is only allocated on demand, page by page, when the memory is actually accessed or modified. This lazy allocation ensures efficiency while allowing ASan to maintain a large shadow memory space without consuming excessive physical resources.

### 4.2.3 Shadow Memory Translation

A given memory address is mapped to a shadow memory address using the formula  $(\text{Address} \gg 3) + \text{Offset}$  (which is equivalent to  $\text{Address} / 8 + \text{Offset}$ ). This means that one byte in shadow memory represents 8 bytes of actual memory.

Because of this 1:8 mapping ratio, the paper allocates one-eighth of the virtual address space for shadow memory—just enough to cover the entire range of possible memory addresses in the program.

The `Offset` represents the starting address of the shadow memory. A simplified interpretation of the mapping is: `SHADOW_MEMORY[Addr >> 3]`.

There's a subtle issue when marking shadow memory. As shown in the figure below, both `arr[0]` and `arr[1]` map to index 0 in the shadow memory. Similarly, `arr[2]` maps to index 2 in the shadow memory. Because of the 1:8 mapping ratio, the next 4 bytes following `arr[2]` (including `arr[3]`) are also covered by the same shadow memory byte at index 2. As a result, OOB access to `arr[3]` may go undetected.

To address this, boundary poisoning is introduced, allowing finer-grained marking of the exact boundaries of valid memory and detecting partial OOB accesses.

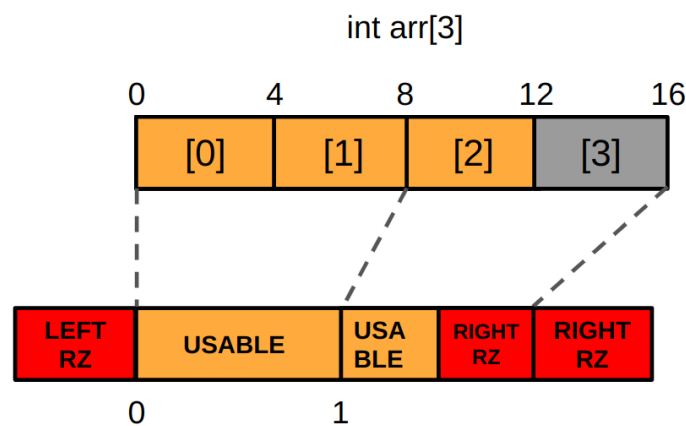


Figure 4.3: Boundary poisoning

The boundary within a shadow memory byte is determined using a bitwise AND operation: `address & 0x07`. For example, if the last address of the variable `arr` is 12, the calculation `12 & 0x07` results in 4. This means that only the first 4 bytes within that 8-byte block are valid. To reflect this, the value `0x04` is written into index 1 of the shadow memory. This indicates that only 4 bytes are part of the usable region. When the variable `arr` is accessed, the instrumented code (explained later) performs a check based on this shadow memory value to determine if the access is valid.

```
1 ShadowAddr = (Addr >> 3) + Offset;
2 k = *ShadowAddr;
3 if (k != 0 && ((Addr & 7) + AccessSize > k))
4     ReportAndCrash(Addr);
```

Consider an access to `arr[3]`. In this case:

- `ShadowAddr = 12 >> 3 = 1`
- `k = 4` (the value stored in shadow memory at index 1, indicating that only the first 4 bytes are valid)
- `AccessSize = 4` (since the array holds integers, each access is 4 bytes)
- At runtime, if condition evaluates true: `if (4 != 0 && (4 + 4 > 4))`
- Correctly triggers an OOB report for `arr[3]`

Now, consider a valid access to `arr[2]`:

- `ShadowAddr = 8 >> 3 = 1`
- `k = 4`
- `AccessSize = 4`
- `if (4 != 0 && (4 + 0 > 4)) → false`
- So the access is considered safe, and no error is reported.

#### Remark 4.1

#### Access range is important

An important observation is that if the `AccessSize` is 5, the condition evaluates to true: `if (4 != 0 && (5 + 0 > 4))`. This correctly triggers an OOB report, since accessing 5 bytes starting at address 8 would reach address 13—beyond the valid region—which is not part of the usable memory.

## 4.2.4 Implementation

ASan is implemented with two components, instrumentation and runtime.

- **Instrumentation Module:** *Instrument runtime call for all memory access operations*
- **Runtime Library:** *Assert if the access is legal*

### 4.2.4.1 ASan Instrumentation

In LLVM IR, memory access operations are represented by two instructions: store and load.

Arrays can be allocated either on the **stack** or on the **heap**. The key difference lies in when their sizes are determined: stack-allocated arrays have a *statically known size* at compile time, whereas heap-allocated arrays are created with a *dynamically determined size* at runtime.

For heap allocations, memory allocation and deallocation are performed through OS-level requests (e.g., `malloc()` / `free()`), which our runtime library intercepts to manage redzones and shadow memory.

In contrast, for statically sized stack arrays, the compiler knows the size at compile time. During instrumentation, the module will insert calls to our runtime to initialize the corresponding shadow memory for the static array.

Here is the main logic of ASan instrumentation module.

```

1  #[derive(Default)]
2  pub struct ASANModule {}
3
4  impl InstrumentModule for ASANModule {
5      fn instrument<'ctx>(<img alt="Rust logo" data-bbox="788 681 805 698"/> Rust
6          &self,
7          context: &'ctx Context,
8          module: &Module<'ctx>,
9          builder: &Builder<'ctx>,
10     ) -> Result<()> {
11         let mut filename_str_ptr = None;
12         let funcs: Vec<_> = module.get_functions().collect();

```



```

13     for func in funcs {
14         // Skip funcs without bodies or those we've added
15         if can_skip_instrument(&func) {
16             continue;
17         }
18         set_filename(module, builder, &mut filename_str_ptr,
19             &func)?;
19         let mut replaced_alloca = HashMap::new();
20         let mut instrumented_blks = HashSet::new();
21         for basic_blk in func.get_basic_blocks() {
22             if instrumented_blks.contains(&basic_blk) {
23                 continue;
24             }
25             for instr in basic_blk.get_instructions() {
26                 match instr.get_opcode() {
27                     // install asan check
28                     InstructionOpcode::Load => {
29                         handle_load(context, module, builder,
30                             filename_str_ptr, &instr)?;
31                     }
32                     InstructionOpcode::Store => {
33                         handle_store(context, module, builder,
34                             filename_str_ptr, &instr)?;
35                     }
36                     InstructionOpcode::Alloca => {
37                         handle_alloca(context, module, builder,
38                             &mut replaced_alloca, &instr)?;
39                     }
40                     // Replace all uses of origin static object
41                     // with newly allocated object's pointer
42                     InstructionOpcode::Call => {
43                         handle_call(context, builder,
44                             &replaced_alloca, &instr)?;
45                     }
46                     InstructionOpcode::GetElementPtr => {
47                         handle_gep(context, builder,
48                             &replaced_alloca, &instr)?;
49                     }
50                     _ => {}
51                 }
52             }
53             instrumented_blks.insert(basic_blk);
54         }
55     }

```

```

49     }
50     // Verify instrumented IRs
51     module_verify(module)
52 }
53 }

```

The instrumentation targets five key LLVM instructions: load, store, alloca, call, and getelementptr. Let's begin by examining the load and store instructions. As their names suggest, both are directly related to memory access operations. We'll start with the `handle_load()` function to understand how load instructions are handled.

```

1  // instrument/src/inkwell_intrinsic.rs 🔒 Rust
2  pub fn get_ptr_operand<'ctx>(instr: &InstructionValue<'ctx>, idx:
   u32) -> PointerValue<'ctx> {
3      instr
4          .get_operand(idx)
5          .unwrap()
6          .left()
7          .unwrap()
8          .into_pointer_value()
9  }
10
11 // instrument/src/asan.rs
12 fn build_memcheck<'ctx>(
13     context: &'ctx Context,
14     module: &Module<'ctx>,
15     builder: &Builder<'ctx>,
16     filename_str_ptr: Option<GlobalValue<'ctx>>,
17     instr: &InstructionValue<'ctx>,
18     ptr: PointerValue<'ctx>,
19     access_size: IntValue<'ctx>,
20 ) -> Result<()> {
21     // the index value must be integer type when using it as array
   index (RHS) value (e.g., int val = arr[idx] + 1)
22     builder.position_before(&instr);
23     build_asan_mem_check(
24         context,
25         module,
26         builder,
27         &filename_str_ptr.unwrap(),
28         ptr,
29         access_size,
30     )?;

```

```

31     Ok(())
32 }
33
34 // instrument/src/asan.rs
35 fn handle_load<'ctx>(<
36     context: &'ctx Context,
37     module: &Module<'ctx>,
38     builder: &Builder<'ctx>,
39     filename_str_ptr: Option<GlobalValue<'ctx>>,
40     instr: &InstructionValue<'ctx>,
41 ) -> Result<()> {
42     let ptr = get_ptr_operand(&instr, 0);
43     if instr.get_type().is_int_type() {
44         let access_size =
45             instr.get_type().into_int_type().size_of();
46         build_memcheck(
47             context,
48             module,
49             builder,
50             filename_str_ptr,
51             &instr,
52             ptr,
53             access_size,
54         )?;
55     }
56     Ok(())
57 }

```

From the load instruction, we can extract the pointer operand that specifies the memory location to read from. This is done using the `get_ptr_operand()` function, which retrieves the operand from the instruction. In a load instruction, the first operand is always the pointer to the data being accessed.

The access size can be determined by querying the type of the loaded value, which is straightforward for statically sized arrays since their types are known at compile time. For simplicity, we handle only integer arrays. However, readers can extend the implementation to support more general types.

Finally, we instrument a call to the runtime to check whether the memory access falls within a valid range. In summary, from a load instruction, we extract both the pointer address and the access size and pass them to the runtime call arguments.

```

1 // instrument/src/inkwell_intrinsic.rs
2 fn get_asan_mem_check_func<'ctx>(<
3     context: &'ctx Context,
4     module: &Module<'ctx>,

```



```

5  ) -> FunctionValue<'ctx> {
6      match get_func(module, ASAN_MEM_CHECK) {
7          Some(func) => func,
8          None => {
9              let asan_mem_check_func_typ =
10                 context.void_type().fn_type(
11                     &[
12                         context.ptr_type(AddressSpace::default()).into(),
13                         context.ptr_type(AddressSpace::default()).into(),
14                         context.i64_type().into(),
15                     ],
16                     false,
17                 );
18             module.add_function(ASAN_MEM_CHECK,
19                                 asan_mem_check_func_typ, None)
20         }
21     }
22 // instrument/src/inkwell_intrinsic.rs
23 pub fn build_asan_mem_check<'ctx>(<
24     context: &'ctx Context,
25     module: &Module<'ctx>,
26     builder: &Builder<'ctx>,
27     filename_str_ptr: &GlobalValue,
28     ptr: PointerValue,
29     access_size: IntValue<'ctx>,
30 ) -> Result<()> {
31     let asan_mem_check = get_asan_mem_check_func(context, module);
32     builder.build_call(
33         asan_mem_check,
34         &[
35             filename_str_ptr.as_pointer_value().into(),
36             ptr.into(),
37             access_size.into(),
38         ],
39         "",
40     )?;
41     Ok(())
42 }

```

Building the runtime call for memory access assertions is straightforward. As with the Coverage module, the first step is to define the signature of the runtime function. This

function takes three parameters: a pointer to the file name, the pointer to the array being accessed, and the size of the access.

Handling the store instruction follows a similar process, with the main difference being how the parameters are retrieved from the instruction.

```

1  // instrument/src/inkwell_intrinsic.rs Rust
2  fn handle_store<'ctx>(
3      context: &'ctx Context,
4      module: &Module<'ctx>,
5      builder: &Builder<'ctx>,
6      filename_str_ptr: Option<GlobalValue<'ctx>>,
7      instr: &InstructionValue<'ctx>,
8  ) -> Result<()> {
9      let offset = instr.get_operand(0).unwrap().left().unwrap();
10     let ptr = get_ptr_operand(&instr, 1);
11
12     // pointer value is used as LHS in `store` instruction. (e.g.,
13     // arr[idx] = value)
14     // any type of access size is valid (i.e., `= value`)
15     let access_size = offset.get_type().size_of().unwrap();
16     build_memcheck(
17         context,
18         module,
19         builder,
20         filename_str_ptr,
21         &instr,
22         ptr,
23         access_size,
24     )?;
25     Ok(())
26 }
```

The remaining instructions to handle are `alloca`, `call`, and `getelementptr`. Let's start by exploring the `alloca` instruction.

A statically allocated array—for example, `int array[3] = {1, 2, 3}`—is represented in LLVM IR as `%1 = alloca [3 x i32], align 4`. To support redzones, our instrumentation module needs to replace this instruction with a new allocation that includes additional space for redzones while preserving the original semantics of the initialization. After instrumentation, the modified instruction becomes `%1 = alloca [19 x i32], align 4`, which allocates a total of 76 bytes—comprising the original 12 bytes and an additional 64 bytes for the redzones.

After performing the replacement, we insert a call to a runtime function to initialize the corresponding shadow memory.

Here's how this can be implemented:

```
1 // instrument/src/inkwell_intrinsic.rs Rust
2 fn handle_alloca<'ctx>(
3     context: &'ctx Context,
4     module: &Module<'ctx>,
5     builder: &Builder<'ctx>,
6     replaced_alloca: &mut HashMap<PointerValue<'ctx>,
7         (PointerValue<'ctx>, IntType<'ctx>)>,
8     instr: &InstructionValue<'ctx>,
9 ) -> Result<()> {
10     // 1. allocate [ redzone | usable | redzone ]
11     let static_arr_kind = instr.get_allocated_type().unwrap();
12     if !static_arr_kind.is_array_type() {
13         return Ok(());
14     }
15     let static_arr = static_arr_kind.into_array_type();
16     let arr_len = static_arr.len();
17     let arr_typ = static_arr.get_element_type();
18     let elem_typ = arr_typ.into_int_type();
19     let elem_typ_byte = elem_typ.get_bit_width() / 8;
20     let rz_size = REDZONE_SIZE / elem_typ_byte;
21
22     builder.position_before(&instr);
23     let new_alloca =
24     builder.build_alloca(elem_typ.array_type(rz_size + arr_len +
25         rz_size), "");
26     let new_alloca_instr = new_alloca.as_instruction().unwrap();
27     new_alloca_instr
28     .set_alignment(instr.get_alignment().unwrap())
29     .unwrap();
30     instr.replace_all_uses_with(&new_alloca_instr);
31     instr.erase_from_basic_block();
32
33     let next_instr_of_new_alloc =
34     new_alloca_instr.get_next_instruction().unwrap();
35     builder.position_before(&next_instr_of_new_alloc);
36
37     // 2. mark redzones and set shadow memory
38     build_asan_init_redzone(
39         context,
40         module,
41         builder,
42         new_alloca,
```

```

39     static_arr.size_of().unwrap(),
40     )?;
41
42     // 3. add redzone size to allocated pointer to correctly set the
43     // usable pointer
44     let new_alloca_ptr =
45     builder.build_alloca(context.ptr_type(AddressSpace::default()),
46     "");
47     let zero_offset = context.i64_type().const_int(0, false);
48     let new_alloca_start_ptr = unsafe {
49     builder.build_in_bounds_gep(
50     elem_typ.array_type(rz_size + arr_len + rz_size),
51     new_alloca,
52     &[zero_offset, zero_offset],
53     "",
54     )?
55     };
56     let rz_offset = context
57     .i64_type()
58     .const_int((REDZONE_SIZE / elem_typ_byte).into(), false);
59     let usable_ptr =
60     unsafe { builder.build_gep(elem_typ, new_alloca_start_ptr,
61     &[rz_offset], "")? };
62     builder.build_store(new_alloca_ptr, usable_ptr)?;
63
64     replaced_alloca.insert(new_alloca, (new_alloca_ptr, elem_typ));
65     Ok(())
66 }

```

First, we determine the type being allocated. If it's not an array type, there's nothing to instrument, so we return early.

If it is an array, we extract relevant information such as the array's length, its overall type, and the element type. These are needed to construct a new `alloca` instruction that accounts for the additional memory required for redzones.

We then create the new `alloca` instruction with the extended size and replace the original allocation uses using `replace_all_uses_with()` and `erase_from_basic_block()`, both provided by the LLVM IR binding library (`inkwell` crate).

Once the replacement is complete, we insert a call to the runtime function to initialize the corresponding shadow memory.

Finally, to preserve the original semantics, we adjust the returned pointer so that it points to the start of the usable memory region—just as the original array would. This logic is implemented in lines 42–57.

Since we replaced the original array with the instrumented one, only the `alloca` variable itself is updated. However, arrays are typically accessed through pointers, not

directly via the `alloca` instruction—because array indexing requires pointer arithmetic. Therefore, we must also update all pointer-based uses of the original array to reference the instrumented version.

This is handled by the `handle_call()` and `handle_gep()` functions.

### Remark 4.2

### Summary

The instrumentation module primarily focuses on three tasks:

1. Replacing the original static array with a newly allocated one that includes additional space for redzones.
2. Updating all uses of the original array to reference the new, instrumented one.
3. Inserting a runtime call to initialize the corresponding shadow memory.

With this, static arrays are fully handled at instrumentation time, while dynamic arrays are managed by the runtime during program execution.

#### 4.2.4.2 ASan Runtime

Runtime tasks are mainly separated into three:

- **Shadow Memory Initialization**
  - Use special markers to identify memory regions: redzones are marked with a redzone marker, and the usable region is initialized with zeros.
- **Memory Access Assertion**
  - Assert if the memory access is legal
- **Allocator Overloading**
  - Override `malloc()`, `realloc()`, `free()`, and C standard library such as `strcpy()`

Let's take a look how shadow memory is initialized:

```

1  // asan_runtime/src/asan_hook.rs
2  #[no_mangle]
3  pub unsafe extern "C" fn __asan_init_redzone(
4      raw_ptr: *mut u8,
5      usable_size: size_t,
6      alloc_kind: u8,
7  ) -> *mut c_void {
8      if raw_ptr.is_null() {
9          return ptr::null_mut();
10     }
11
12     let usable_ptr = unsafe { raw_ptr.add(REDZONE_SIZE) };
13
14     // in a right manner of implementation, this memset seems
15     // redundant
16     // without this unused explicit memory initialization, it may
17     // trigger segfault

```



```

16 // you may try comment out these two lines and use z3 library,
   // then segfault will be triggered
17 // when `free` system call is invoked.
18 // exact reasoning has not been found yet
19 // also, initialization with zero value only works empirically
20 libc::memset(raw_ptr as *mut c_void, 0, REDZONE_SIZE);
21 libc::memset(raw_ptr.add(REDZONE_SIZE + usable_size) as *mut
   c_void, 0, REDZONE_SIZE);
22
23 // initialize shadow memory
24 poison_shadow_allocated(raw_ptr as usize, usable_size,
   alloc_kind);
25 // record allocated size
26 ALLOC_MAP
27     .lock()
28     .unwrap()
29     .insert(usable_ptr as usize, usable_size)
30     .expect(ALLOC_MAP_INSERT_ERR_STR);
31
32 // return usable region pointer
33 usable_ptr as *mut c_void
34 }
35
36 pub fn convert_to_shadow_idx(addr: usize) -> usize {
37     addr >> SHADOW_SCALE
38 }
39
40 fn write_shadow_mem(shadow_mem: &mut [i8], idx: usize, val: i8) {
41     shadow_mem[idx % SHADOW_SIZE] = val;
42 }

```

Redzone memory is initially filled with zeros. The `poison_shadow_allocated()` function is then used to initialize the corresponding shadow memory.

```

1 // asan_runtime/src/asan_hook.rs Rust
2 lazy_static::lazy_static! {
3     pub static ref SHADOW_MEMORY: Mutex<Vec<i8>> = Mutex::new(vec!
   [0i8; SHADOW_SIZE]);
4     pub static ref ALLOC_MAP: Mutex<FnvIndexMap<usize, usize,
   ALLOC_MAP_SIZE>> = Mutex::new(FnvIndexMap::new());
5 }
6
7 /// Poisons left and right redzones and make clean for usable region

```

```

8  fn poison_shadow_allocated(raw_ptr: usize, usable_size: usize,
   alloc_kind: u8) {
9      //
   |-----|-----|-----|
10     // ^      LEFT_RZ      USABLE      ^
   RIGHT_RZ      ^
11     //
12     // precisely allocate shadow byte for each boundary(^)
13
14     let usable_ptr = raw_ptr + REDZONE_SIZE;
15
16     let (left_rz_marker, right_rz_marker) = match alloc_kind {
17         1 => (STACK_LEFT_REDZONE_MARKER,
18             STACK_RIGHT_REDZONE_MARKER),
19         2 => (HEAP_LEFT_REDZONE_MARKER, HEAP_RIGHT_REDZONE_MARKER),
20         _ => panic!("unknown alloc kind"),
21     };
22
23     // 1. Poison left redzone
24     let left_start = raw_ptr;
25     let left_end = usable_ptr;
26     let shadow_left_start = convert_to_shadow_idx(left_start);
27     let shadow_left_end = convert_to_shadow_idx(left_end);
28     let mut shadow_mem = SHADOW_MEMORY.lock().unwrap();
29
30     set_boundary_poison_byte(
31         &mut shadow_mem,
32         left_start,
33         shadow_left_start,
34         left_rz_marker,
35     );
36     for i in (shadow_left_start + 1)..shadow_left_end {
37         write_shadow_mem(&mut shadow_mem, i, left_rz_marker);
38     }
39
40     // 2. Unpoison usable region
41     let usable_start = usable_ptr;
42     let usable_end = usable_ptr + usable_size;
43     let shadow_usable_start = convert_to_shadow_idx(usable_start);
44     let shadow_usable_end = convert_to_shadow_idx(usable_end);
45
46     for i in shadow_usable_start..shadow_usable_end {

```

```
46     write_shadow_mem(&mut shadow_mem, i, CLEAN_BYTE_MARKER);
47 }
48 set_bounary_poison_byte(
49     &mut shadow_mem,
50     usable_end,
51     shadow_usable_end,
52     right_rz_marker,
53 );
54
55 // 3. Poison right redzone
56 let right_start = usable_end;
57 let right_end = right_start + REDZONE_SIZE;
58 let shadow_right_start = convert_to_shadow_idx(right_start);
59 let shadow_right_end = convert_to_shadow_idx(right_end);
60
61 for i in (shadow_right_start + 1)..shadow_right_end {
62     write_shadow_mem(&mut shadow_mem, i, right_rz_marker);
63 }
64 set_bounary_poison_byte(
65     &mut shadow_mem,
66     right_end,
67     shadow_right_end,
68     right_rz_marker,
69 );
70 }
71
72 // asan_runtime/src/asan_hook.rs
73 fn set_bounary_poison_byte(
74     shadow_mem: &mut [i8],
75     start: usize,
76     shadow_start: usize,
77     rz_marker: i8,
78 ) {
79     let remaining = start & 0x07;
80     if remaining != 0 {
81         write_shadow_mem(shadow_mem, shadow_start, remaining as i8);
82     } else {
83         write_shadow_mem(shadow_mem, shadow_start, rz_marker);
84     }
85 }
86
```

```

87 // asan_runtime/src/asan_hook.rs
88 pub fn convert_to_shadow_idx(addr: usize) -> usize {
89     addr >> SHADOW_SCALE
90 }
91
92 // asan_runtime/src/asan_hook.rs|
93 fn write_shadow_mem(shadow_mem: &mut [i8], idx: usize, val: i8) {
94     shadow_mem[idx % SHADOW_SIZE] = val;
95 }

```

- L17 - 32 fills LEFT\_REDZONE marker in left redzone area
- L34 - 48 fills CLEAN\_BYTE marker in usable region
- L50 - 64 fills RIGHT\_REDZONE marker in right redzone area

A key point to note is that, as mentioned earlier, precise boundary poisoning is crucial to accurately define the boundaries of each memory region.

ALLOC\_MAP tracks usable region's pointer and its size where is necessary in `realloc()` and `free()` function.

Up to this point, our ASan runtime has only handled static arrays. Now it's time to extend support to dynamic arrays. Since heap memory is managed through standard C functions like `malloc()`, `realloc()`, and `free()`, we need to intercept and override these calls. The first target for interception is `malloc()`.

```

1 // asan_runtime/src/asan_intrinsic.rs Rust
2
3 type MallocFn = unsafe extern "C" fn(size_t) -> *mut c_void;
4 type ReallocFn = unsafe extern "C" fn(ptr: *mut c_void, size:
  size_t) -> *mut c_void;
5 type FreeFn = unsafe extern "C" fn(ptr: *mut c_void);
6 type StrcpyFn = unsafe extern "C" fn(dest: *mut c_char, src: *const
  c_char) -> *mut c_char;
7
8 lazy_static::lazy_static! {
9     static ref ALLOC_MAP: Mutex<HashMap<usize, usize>> =
      Mutex::new(HashMap::new());
10    static ref MALLOC: Mutex<Option<MallocFn>> = Mutex::new(None);
11    static ref REALLOC: Mutex<Option<ReallocFn>> = Mutex::new(None);
12    static ref FREE: Mutex<Option<FreeFn>> = Mutex::new(None);
13    static ref STRCPY: Mutex<Option<StrcpyFn>> = Mutex::new(None);
14 }
15
16 pub fn get_cmalloc() -> MallocFn {
17     let mut cmalloc = MALLOC.lock().unwrap();
18     if cmalloc.is_none() {
19         let buf = get_malloc_cstr();

```

```

20     let malloc_sym = unsafe { dlsym(RTLD_NEXT, buf.as_ptr()) };
21     let malloc_fn =
22         unsafe { std::mem::transmute::<*const (),
23             MallocFn>(malloc_sym as *const ()) };
24     *cmalloc = Some(malloc_fn);
25 }
26 cmalloc.unwrap()
27 }
28 // asan_runtime/src/asan_hook.rs
29
30 thread_local! {
31     pub static MALLOC_REENTERED: Mutex<bool> = const
32     { Mutex::new(false) }
33 }
34 #[no_mangle]
35 pub unsafe extern "C" fn malloc(size: size_t) -> *mut c_void {
36     let cmalloc = get_cmalloc();
37
38     MALLOC_REENTERED.with(|re_enter| {
39         let is_reenter = re_enter.lock().unwrap();
40         if *is_reenter {
41             cmalloc(size)
42         } else {
43             let total_size = size + REDZONE_SIZE * 2;
44             let raw_ptr = cmalloc(total_size) as *mut u8;
45             __asan_init_redzone(raw_ptr, size, ALLOC_HEAP)
46         }
47     })
48 }

```

The main responsibility of the overridden `malloc()` function is to allocate additional memory to account for both left and right redzones, and then initialize the corresponding shadow memory.

The standard memory allocation process still applies. Therefore, we must request a larger allocation from the operating system. To do this, the first step is to retrieve the original `malloc()` function dynamically from the symbol table.

Once the original `malloc()` is obtained, we request a total allocation of `size + 2 * REDZONE_SIZE`, and then call the shadow memory initialization routine with the allocation kind `ALLOC_HEAP` which indicates whether the memory is stack or heap.

For other standard functions, the same process is repeated per function.

```
1 // asan_runtime/src/asan_intrinsic.rs
```

 Rust

```
2
3 pub fn get_malloc_cstr() -> [i8; 7] {
4     let fn_str = b"malloc";
5     let mut buf = [0 as c_char; 7];
6     for (i, &b) in fn_str.iter().enumerate() {
7         buf[i] = b as c_char;
8     }
9     buf
10 }
11
12 pub fn get_realloc_str() -> [i8; 8] {
13     let fn_str = b"realloc";
14     let mut buf = [0 as c_char; 8];
15     for (i, &b) in fn_str.iter().enumerate() {
16         buf[i] = b as c_char;
17     }
18     buf
19 }
20
21 pub fn get_free_str() -> [i8; 5] {
22     let fn_str = b"free";
23     let mut buf = [0 as c_char; 5];
24     for (i, &b) in fn_str.iter().enumerate() {
25         buf[i] = b as c_char;
26     }
27     buf
28 }
29
30 pub fn get_strcpy_str() -> [i8; 7] {
31     let fn_str = b"strcpy";
32     let mut buf = [0 as c_char; 7];
33     for (i, &b) in fn_str.iter().enumerate() {
34         buf[i] = b as c_char;
35     }
36     buf
37 }
38
39 pub fn get_cmalloc() -> MallocFn {
40     let mut cmalloc = MALLOC.lock().unwrap();
41     if cmalloc.is_none() {
42         let buf = get_malloc_cstr();
```

```

43     let malloc_sym = unsafe { dlsym(RTLD_NEXT, buf.as_ptr()) };
44     let malloc_fn =
45         unsafe { std::mem::transmute::<*const (),
46             MallocFn>(malloc_sym as *const ()) };
47     *cmalloc = Some(malloc_fn);
48 }
49 cmalloc.unwrap()
50
51 pub fn get_crealloc() -> ReallocFn {
52     let mut crealloc = REALLOC.lock().unwrap();
53     if crealloc.is_none() {
54         let buf = get_realloc_str();
55         let realloc_sym = unsafe { dlsym(RTLD_NEXT, buf.as_ptr()) };
56         let realloc_fn =
57             unsafe { std::mem::transmute::<*const (),
58                 ReallocFn>(realloc_sym as *const ()) };
59         *crealloc = Some(realloc_fn);
60     }
61     crealloc.unwrap()
62 }
63
64 pub fn get_cfree() -> FreeFn {
65     let mut cfree = FREE.lock().unwrap();
66     if cfree.is_none() {
67         let buf = get_free_str();
68         let free_sym = unsafe { dlsym(RTLD_NEXT, buf.as_ptr()) };
69         let free_fn = unsafe { std::mem::transmute::<*const (),
70             FreeFn>(free_sym as *const ()) };
71         *cfree = Some(free_fn);
72     }
73     cfree.unwrap()
74 }
75
76 pub fn get_strcpy() -> StrcpyFn {
77     let mut strcpy = STRCPY.lock().unwrap();
78     if strcpy.is_none() {
79         let buf = get_strcpy_str();
80         let strcpy_sym = unsafe { dlsym(RTLD_NEXT, buf.as_ptr()) };
81         let strcpy_fn =
82             unsafe { std::mem::transmute::<*const (),
83                 StrcpyFn>(strcpy_sym as *const ()) };

```

```

81     *cstrcpy = Some(strcpy_fn);
82 }
83 cstrcpy.unwrap()
84 }

```

free() implementation is as follows:

```

1  // asan_runtime/src/asan_hook.rs Rust
2  #[no_mangle]
3  pub unsafe extern "C" fn free(ptr: *mut c_void) {
4      if ptr.is_null() {
5          return;
6      }
7      let cfree = get_cfree();
8      let mut alloc_map = ALLOC_MAP.lock().unwrap();
9      if let Some(size) = alloc_map.remove(&(ptr as usize)) {
10         poison_shadow_freed(ptr, size);
11         cfree(ptr.sub(REDZONE_SIZE));
12     } else {
13         cfree(ptr);
14     }
15 }

```

The first step is to obtain the standard free() symbol. Since our custom malloc() handler records each allocation in the ALLOC\_MAP structure—including the usable pointer and its size—we can track all memory managed by our runtime.

When free() is called, we check whether the pointer was allocated through our malloc() handler by querying ALLOC\_MAP. If it was, we retrieve the allocation size and mark the entire region—including the usable portion and both redzones—with a FREED marker.

Any subsequent access to this freed region will be detected by our access checker, allowing us to report UAF. The process of marking shadow memory as FREED is as follows:

```

1  // asan_runtime/src/asan_hook.rs Rust
2  /// Make all regions as poisoned
3  fn poison_shadow_freed(usable_ptr: *mut c_void, size: usize) {
4      let start = unsafe { usable_ptr.sub(REDZONE_SIZE) };
5      let shadow_start = convert_to_shadow_idx(start as usize);
6      let shadow_end =
7          unsafe { convert_to_shadow_idx(start.add(REDZONE_SIZE + size
8              + REDZONE_SIZE) as usize) };
9      let mut shadow = SHADOW_MEMORY.lock().unwrap();
10     set_boundary_poison_byte(&mut shadow, start as usize,
11         shadow_start, FREED_MARKER);

```



```

10     for i in (shadow_start + 1)..shadow_end {
11         write_shadow_mem(&mut shadow, i, FREED_MARKER);
12     }
13     set_boundary_poison_byte(&mut shadow, start as usize, shadow_end,
14                             FREED_MARKER);
15 }

```

The third function we intercept is `realloc()`. The behavior of `realloc()` can be broken down into two main tasks:

1. Allocate a new memory block of the requested size and copy the original data into it.
2. Free the original memory block.

Our implementation must preserve this behavior while adding shadow memory handling. Specifically, we invoke our custom `malloc()` to perform the new allocation and initialize the corresponding shadow memory. Then, we call our custom `free()` handler to mark the original memory (including its redzones) as `FREED` in shadow memory.

```

1  // asan_runtime/src/asan_hook.rs Rust
2  #[no_mangle]
3  pub unsafe extern "C" fn realloc(ptr: *mut c_void, size: size_t) ->
4      *mut c_void {
5      if ptr.is_null() {
6          return malloc(size);
7      }
8      let mut alloc_map = ALLOC_MAP.lock().unwrap();
9      if let Some(org_size) = alloc_map.remove(&(ptr as usize)) {
10         // drop locked variables because the following code will
11         // require `free` call
12         drop(alloc_map);
13         poison_shadow_freed(ptr, size);
14         let usable_ptr = malloc(size);
15         ptr::copy(ptr, usable_ptr, org_size);
16         free(ptr.sub(REDZONE_SIZE));
17         usable_ptr
18     } else {
19         let crealloc = get_crealloc();
20         crealloc(ptr, size)
21     }
22 }

```

Now, our implementation til now handles both static and dynamic array.

```

1  #include <stdlib.h> C
2
3  int main() {

```

```

4     int* arr = (int*)malloc(sizeof(int) * 3);
5     arr[0] = 1;
6     arr[1] = 1;
7     arr[2] = 1;
8     arr[3] = 1; // 00B report
9 }

```

OOB report at line number 8:

```

1  > LD_LIBRARY_PATH=./ ./a.out
2  [ASAN] invalid memory access detected at kkk.c: 0x640c601292cc
3      0: asan_runtime::asan_runtime::report_asan_violated::{closure}}
4          at .../asan_runtime/src/asan_runtime.rs:46:18
5      1: std::thread::local::LocalKey<T>::try_with
6          at /build/rustc-1.84.1-src/library/std/src/thread/
          local.rs:283:12
7      2: std::thread::local::LocalKey<T>::with
8          at /build/rustc-1.84.1-src/library/std/src/thread/
          local.rs:260:9
9      3: asan_runtime::asan_runtime::report_asan_violated
10         at .../asan_runtime/src/asan_runtime.rs:44:5
11      4: __asan_mem_check
12         at .../asan_runtime/src/asan_runtime.rs:30:9
13      5: main
14         at ./asan-example2.c:8:12
15      6: __libc_start_call_main
16      7: __libc_start_main_alias_2
17      8: _start

```

Now the memory range assertion is introduced:

```

1  // asan_runtime/src/asan_runtime.rs  Rust
2  #[no_mangle]
3  pub extern "C" fn __asan_mem_check(file_ptr: *const libc::c_char,
4      addr: usize, access_size: usize) {
5      if file_ptr.is_null() || addr == 0 {
6          return;
7      }
8      let filename = cstr_to_string(file_ptr);
9      let shadow_idx = convert_to_shadow_idx(addr);
10     let shadow_val = SHADOW_MEMORY.lock().unwrap()[shadow_idx %
        SHADOW_SIZE];
11     if shadow_val != CLEAN_BYTE_MARKER && ((addr & 0x07) +
        access_size) as i8 > shadow_val {

```

```
12     report_asan_violated(&filename, addr);
13 }
14 }
15
16 fn report_asan_violated(filename: &str, addr: usize) {
17     if is_test_enabled() {
18         eprintln!("[ASAN] invalid memory access detected at {}",
19             filename);
20     } else {
21         eprintln!(
22             "[ASAN] invalid memory access detected at {}: 0x{:x}",
23             filename, addr
24         );
25         // print backtrace
26         MALLOC_REENTERED.with(|re_enter| {
27             *re_enter.lock().unwrap() = true;
28             let bt = Backtrace::force_capture();
29             if is_test_enabled() {
30                 eprintln!("{}", trim_runtime_bt(bt.to_string()));
31             } else {
32                 eprintln!("{bt}");
33             }
34             *re_enter.lock().unwrap() = false;
35         });
36         if !is_test_enabled() {
37             unsafe {
38                 libc::_exit(EXIT_CODE);
39             }
40         }
41     }
42
43 fn is_test_enabled() -> bool {
44     matches!(env::var(ASAN_TEST_ENABLED), Ok(val) if val == "1")
45 }
46
47 fn trim_runtime_bt(bt: String) -> String {
48     bt.lines()
49         .skip_while(|line| !line.contains("__asan_mem_check"))
50         .skip(2)
51         .collect::<Vec<_>>()
```

```
52         .join("\n")
53     }
```

An important detail to note is the use of `MALLOC_REENTERED`. The function `poison_shadow_allocated()` may itself trigger memory allocation when initializing shadow memory. Although shadow memory is declared globally, it is not fully allocated up front—memory is committed on-demand, page by page.

When an OOB violation is reported, functions like `eprintln!` may also perform dynamic memory allocation. This leads to a call to our overridden `malloc()`, which can result in a recursive loop and ultimately cause a stack overflow.

To prevent this, we use `MALLOC_REENTERED` to detect and break the recursion. During reporting, we bypass the normal path and allow a minimal, direct call of standard `malloc()`.

This does not compromise shadow memory management, as the allocation is unrelated to the target program's memory space and occurs just before the program terminates due to the reported error.

Lastly, an important point to note: since each redzone is 32 bytes, totaling 64 bytes, any access beyond the redzone boundaries will go undetected by the runtime. In other words, if OOB access occurs outside the redzones, no report will be triggered.

You can test this behavior by modifying line 8 as follows:

```
1     arr[50] = 1; // Not reported
```

However, there are certain edge cases that our implementation cannot currently detect. For example, consider the following scenario:

```
1 #include <string.h>
2
3 int main() {
4     char buffer[10];
5     strcpy(buffer, "aaaaaaaaa");
6     return 0;
7 }
```

This is clearly an OOB bug. However, the standard `strcpy()` function is part of the C runtime and executes outside our instrumented code, making it invisible to our current instrumentation. Similar issues can arise when other standard library functions perform memory operations outside the scope of our target program code.

To address this, we can intercept these C runtime functions—just as we did with `malloc()`, `free()`, and `realloc()`—by overriding them. The following is our custom implementation of `strcpy()`.

```
1 #[no_mangle]
2 pub unsafe extern "C" fn strcpy(dest: *mut c_char, src: *const
   c_char) -> *mut c_char {
3     let cstrcpy = get_strcpy();
4     let src_len = CStr::from_ptr(src).to_bytes().len();
5     let temp_filename_ptr = c"libc::strcpy".as_ptr() as *const
   c_char;
```

```
6     for i in 0..=src_len {
7         __asan_mem_check(temp_filename_ptr, dest as usize + i, 1);
8     }
9     strcpy(dest, src)
10 }
```

The core behavior of `strcpy()` involves copying elements one by one by iterating through the source array. We insert memory access checks for each index of the destination buffer (as shown in lines 6–8). After verifying the safety of all accesses, we delegate the actual copying to the original `strcpy()` function provided by the C runtime.

#### 4.2.4.3 Summary

We implemented OOB detector using both compile-time instrumentation and a runtime library. The instrumentation module inserts runtime function calls that assert the validity of memory accesses, while the runtime intercepts memory allocation functions and manages shadow memory.

You can refer to the full implementation here:

- Instrumentation module: `instrument/src`
- Runtime library: `asan_runtime/src`

Test cases are available in the `coverage/tools/tests/inputs/asan` directory

#### 4.2.4.4 Homework

**Exercise 4.1** We suggest the following topics for further study and exploration:

- **Topic #1: Global object handling**
  - Our implementation supports both static and heap-allocated arrays. For static arrays, we currently handle only local static arrays, not global ones. However, support for global arrays can be added by extending the instrumentation at the frontend level. We encourage readers to explore this as a possible enhancement.
- **Topic #2: C standard library support**
  - This book demonstrates how to override the C standard library function `strcpy()` as an example. Readers are encouraged to extend this approach to other functions such as `strlen()`, `strcmp()`, and beyond.



## Part IV

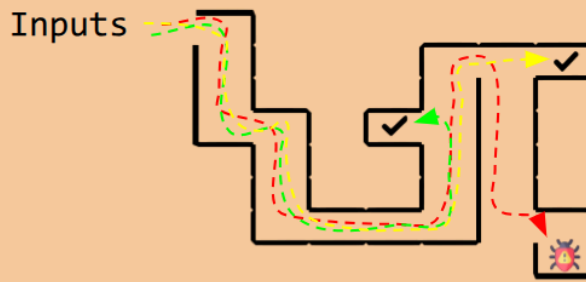
# Fuzzing

<b>5</b>	<b>Fuzzer . . . . .</b>	<b>81</b>
5.1	Introduction to Fuzzer . . . . .	81
5.2	A Big Picture . . . . .	83
5.3	Design and Architecture . . . . .	84
5.4	Implementation . . . . .	87





# “Automating the process of finding a program’s vulnerabilities”



## 5. Fuzzer

### 5.1 Introduction to Fuzzer

Fuzzing [21] is one of the most effective software testing methods and a well-known approach for discovering bugs. It is widely used in both academia and industry.

A naive fuzzer works by generating randomized inputs and feeding them to the target program until a crash is detected.

Let’s imagine what fuzzing would be. The example code below takes user input via stdin and evaluate the value. If the value is 10, then prgoram will crash, otherwise, print the value to terminal. Assume that we don’t know the crash condition (i.e, when `input()` evaluates to 10)

```
1 v = input()
2 if (v == 10) {
3     crash()
4 } else {
5     printf("%d", v);
6 }
```

An example program waits for input. A fuzzer generates random input and feeds it to the program through stdin. The program, previously blocked waiting for input, then proceeds to execute the next portion of its code. Suppose the input is 1; in this case, the program exits normally, indicating no crash. The fuzzer then moves to the next round—say with input 2—and continues generating and feeding inputs. This process is repeated until a crash is triggered.

This represents the original idea behind fuzzing when it was first introduced. Since then, fuzzers have evolved significantly to find crashes more effectively. The figure below illustrates the basic concept of a naive fuzzer versus a modern, off-the-shelf fuzzer.

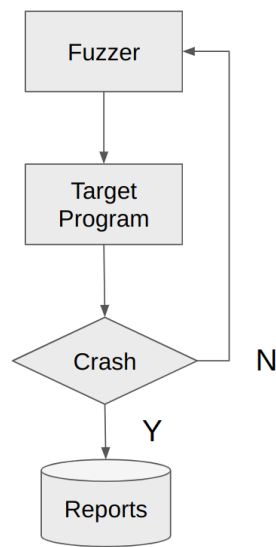


Figure 5.1: A naive fuzzer

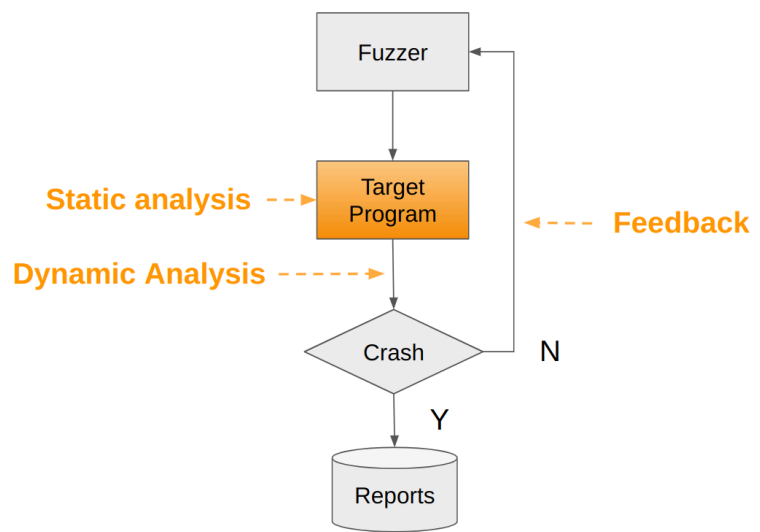


Figure 5.2: Advanced fuzzer

The left figure illustrates an imaginary naive fuzzer, while the right figure shows an off-the-shelf fuzzer that leverages both static and dynamic analysis to understand program structure. These advanced techniques help guide the fuzzing process and significantly improve the effectiveness of the fuzzing campaign. Modern fuzzers contribute not only in generating inputs but also in analyzing programs to extract useful information—such as probability-based heuristics, algorithmic insights, or theoretical models—to make fuzzing more intelligent and targeted.

Fuzzers can be broadly categorized into three types: black-box, white-box, and gray-box.

- A **black-box fuzzer**, as the name suggests, treats the program as an opaque system. It only observes output behavior in response to given inputs.
- A **white-box fuzzer** theoretically has access to all internal details of the program. It can leverage static and dynamic analysis to understand and exploit the program's behavior for optimal input generation.
- A **gray-box fuzzer** operates in between, utilizing partial information—most commonly, code coverage—to guide input generation.

In this book, we will implement a **gray-box fuzzer**.

There are several factors that contribute to an effective fuzzer, with coverage being one of the most widely used and important metrics.

**American Fuzzy Lop (AFL)** [22] is one of the most popular coverage-based fuzzers. It instruments the target program to collect coverage information at runtime, which the fuzzer then uses to guide its exploration toward more interesting code paths.

It's important to note that *coverage-based* does not imply that the fuzzer actively increases code coverage. Instead, **it tracks which parts of the code have been executed and identifies whether new code paths have been discovered. Inputs that reach previously unseen or rarely executed code are assigned higher scores.** In essence, a coverage-based fuzzer determines its fuzzing strategy based on code coverage feedback.

In this book, we use AFL as our reference model and implement an AFL-like fuzzer in Rust. As in previous chapters, our fuzzer consists of two main components: instrumentation and a runtime library. Let's begin by exploring the design of our fuzzer.

## 5.2 A Big Picture

The figure below illustrates the fuzzer we will implement throughout this book. Symbolic execution and crash minimization will be introduced in Chapters 4 and 5, respectively.

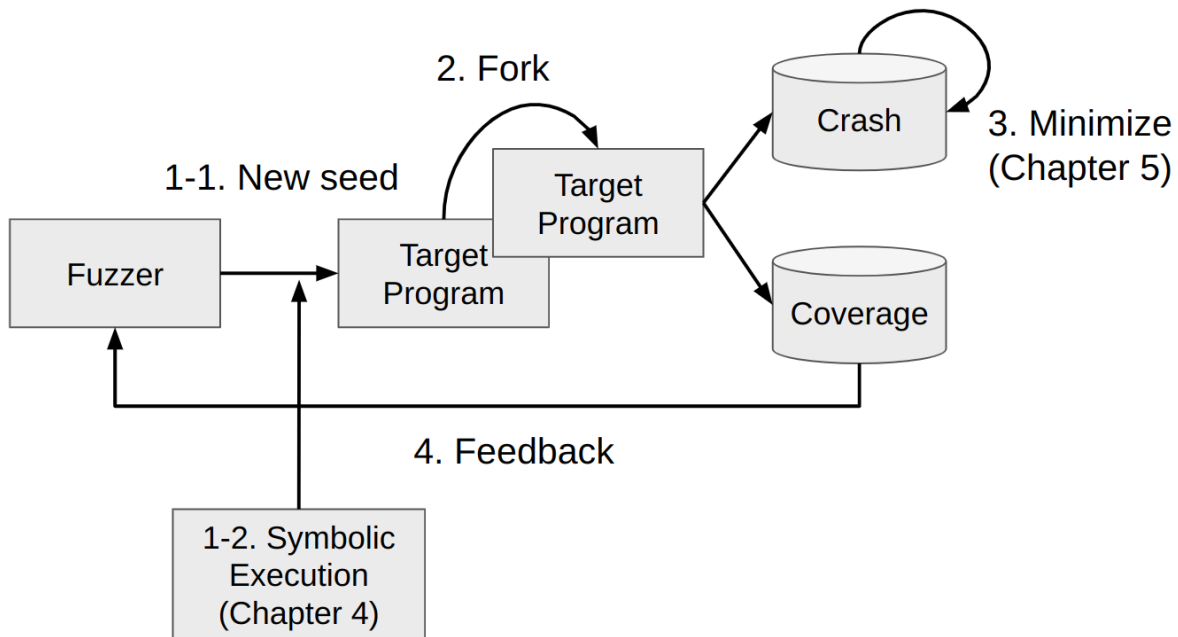


Figure 5.3: A workflow of fuzzing campaign

The fuzzer is an independent binary that operates separately from the instrumentation, runtime, and target program. It generates a new seed input, which is sent to the target program. The target program remains idle until it receives this input from the fuzzer. Upon receiving the seed, it forks itself, and the child process executes using the provided input.

After the runtime is invoked, the execution can result in two outcomes: a crash or coverage information. If the target program crashes, the input is considered failure-inducing. Before saving this input to the file system, a minimizer attempts to reduce it to the smallest form that still triggers the crash, making it easier for developers to analyze the root cause.

The second type of output is the coverage footprint, which indicates which parts of the code were newly explored, including the specific edges visited. The fuzzer uses this information to guide the next iteration by generating a new seed that aims to expand code coverage.

This sequence of generating input, executing the target, collecting output, and guiding further mutations constitutes a single fuzzing iteration, which is then repeated continuously. In addition, the target program's input can also be generated by a symbolic execution runtime, which will be introduced in Chapter 4.

## 5.3 Design and Architecture

### 5.3.1 Branch(Edge) Coverage

The first important topic is branch (or edge) coverage. As explained in this [LLVM-IR tutorial](#), all instructions within a basic block in LLVM IR are guaranteed not to alter control flow. In AFL, the basic block serves as the fundamental unit of coverage. Consider the following example:

	LLVM-IR
1	bb1:
2	...
3	br i1 %x, label %bb2, label %bb3
4	bb2:
5	...
6	bb3:
7	...

There are three basic blocks, and the possible coverage paths are either bb1 -> bb2 or bb1 -> bb3. Coverage is recorded by instrumenting each basic block using the instrumentation module. This coverage data is stored in shared memory, which is accessible to both the fuzzer and its runtime. Each basic block is instrumented as follows:

```
1 cur_location = <COMPILE_TIME_RANDOM>;
2 shared_mem[cur_location ^ prev_location]++;
3 prev_location = cur_location >> 1;
```

The `cur_location` is a randomly generated value assigned at compile time. `shared_mem` refers to the shared memory region used by both the fuzzer and the runtime to track coverage. The index of the current edge is calculated using the XOR operation: `cur_location ^ prev_location`. This XORed value uniquely identifies the edge in the control flow. For example, the edges bb1 -> bb3 and bb2 -> bb3 produce different XOR values and are therefore treated as distinct. Let's deep dive into the semantic.

Basic block	Current location	Prev location	Shared memory index
A	r1	0	$r1 \wedge 0 = r1$
B	r2	$r1 \gg 1$	$r2 \wedge (r1 \gg 1)$
C	r3	$r2 \gg 1$	$r3 \wedge (r2 \gg 1)$

The table assumes the execution flow A -> B -> C. Each basic block is first assigned a unique random identifier at compile time. To track which path (or edge) has been executed, the runtime calculates `cur_location ^ prev_location`, effectively representing the transition from one basic block to another as a unique tuple. For example, in the path A -> B, only one such tuple (AB) is recorded. A second execution of A -> B will compute the same XOR value and thus update the same index in the shared memory.

```
1 `A -> B -> A -> B`
```

First, A -> B creates a new index of the shared memory. B -> A also hits a new index. The remaining one is only A -> B again. However, at this moment, the calculation creates

an index before we seen when we create a first  $A \rightarrow B$ . Then it increases the count of tuple AB. Thus, the shared memory is imagine as follows:

First, the transition  $A \rightarrow B$  generates a new index in the shared memory. Later,  $B \rightarrow A$  results in a different index, representing a distinct edge. When  $A \rightarrow B$  occurs again, the XOR calculation produces the same index as the initial  $A \rightarrow B$  transition. Instead of creating a new entry, it increments the counter for the existing AB tuple. As a result, the shared memory can be visualized as follows:

```
1 shared_mem[index1] = 2 // A -> B
2 shared_mem[index2] = 1 // B -> A
```

Thus, the second  $A \rightarrow B$  transition does not represent new coverage—it is not a new discovery. Finally, let's look at another example provided in the AFL documentation.

```
1 #1: A -> B -> C -> D -> E
2 #2: A -> B -> C -> A -> E
```

Assume that two paths have already been discovered, and the current path does not introduce any new coverage.

```
1 #3: A -> B -> C -> A -> B -> C -> A -> B -> C -> D -> E
```

Let's break down the path into tuples, one step at a time:

- A -> B (#1)
- B -> C (#1)
- C -> D (#1)
- D -> E (#1)
- C -> A (#2)
- A -> E (#2)

We can easily identify that step #3 does not represent new coverage.

### 5.3.2 Bucketing

As explained earlier, coverage information determines the direction of fuzzing. Coverage-based fuzzers primarily prioritize two things: (1) discovering new paths and (2) exploring less frequently executed paths.

A new path is valuable because it represents code that has never been executed before, making it a prime candidate for uncovering potential bugs. For the second case, consider the probability of finding a bug: if you have executed code path A many times but rarely or never executed path B, it makes sense to focus on path B since it is less explored and may hide undiscovered issues.

In summary, discovering new paths is the highest priority. If no new paths are found, the fuzzer then prioritizes paths that have been executed less frequently.

To represent how many times a path is executed, a coarse-grained bucketing scheme is introduced. The bucket consists of 8 elements representing execution count ranges: 1, 2, 3, 4-7, 8-15, 16-31, 32-127, and 128+. Each bucket index is assigned a score from 8 (for the lowest count) down to 1 (for the highest count). For example, if a path is executed only once, it receives a score of 8; if it is executed more than 128 times, it receives a score of 1.

When new coverage is discovered, it is awarded a score higher than 8 to prioritize novel paths.

This process of grouping execution counts into discrete ranges and assigning scores is called bucketing.

### 5.3.3 Fork Server

When the fuzzer generates a new seed, it is fed to the target program. However, if the target program is reloaded as a new process for each seed, a significant amount of time is wasted on process startup overhead. In general, the following tasks are involved each time a new process is loaded:

- Disk I/O: Read the binary from disk and loads into memory
- Initialization:
  - Global variable / Heap / Stack
  - Linker / Loader
  - Resource configuration (file descriptor / socket / etc)

To avoid the repeated overhead of loading a new process for every seed, the fork-server mechanism is introduced.

There are various ways to implement a fork-server. In this book, it is implemented within the runtime. The runtime initially loads the target program into memory once and never reloads it afterward. After loading, the process remains blocked, waiting for a signal—let's call it *wakeup*. Upon receiving this signal, the process forks itself, resulting in two processes: a parent and a child.

The parent process remains blocked until the child process finishes execution. Meanwhile, the child process runs the main logic of the target program using the provided seed. Once the child exits, the parent collects the exit status and signals it back to the fuzzer. This cycle repeats for each seed.

By doing this, the fork-server eliminates the costly process initialization step. Instead of starting a fresh process each time, it efficiently forks from a pre-initialized state. Communication between the fuzzer and the fork-server can be handled via a dedicated shared memory region, which will be explained in detail in the implementation section.

### 5.3.4 Harness

*Think about any program you want to fuzz—and how to feed fuzzing input into it.*

For a target program to accept input from an external source like a fuzzer, it must expose an interface for receiving that input. However, most programs are not designed to cooperate with fuzzing by default. Therefore, we often need to modify or wrap the original target program to enable fuzzing interaction. This wrapper or interface is called a ***harness***.

A harness must be tailored to each specific program. For example, a GUI program may require simulating user interactions such as mouse clicks or keystrokes, while an IoT program might need to simulate external sensor data.

In general, for text user interface (TUI) programs, two common methods are used to deliver input: via `stdin` or via files.

For `stdin`, the harness directs the fuzzing input to the program's standard input file descriptor. The program reads from `stdin` and processes the input through its main logic.

For file-based input, the harness writes the fuzzing input to a file, and the target program reads from that file as if it were regular input. The file content is updated for each fuzzing iteration.

In this book, we focus on two types of harnesses: `stdin`-based and file-based.

### 5.3.5 Oracle

*If a program crashes, how do we determine whether it completed normally or actually crashed?*

To answer this, we need a clear definition of what constitutes a crash. In many security papers, researchers define an oracle—a mechanism for identifying abnormal or interesting behaviors in the target program.

A common oracle is a segmentation fault. If a segmentation fault occurs, we can confidently say the program crashed. However, the absence of a segmentation fault does not necessarily mean the program executed correctly. For instance, in Chapter 2, we discussed OOB bugs: one of the bug category such as segmentation fault. Unfortunately, in C programs, certain OOB accesses may not immediately trigger visible crashes or errors.

To detect such hidden bugs, the target program can be compiled with sanitizers like ASan. With ASan enabled, the program will crash explicitly when an OOB access occurs. The fuzzer can then detect this by checking the program's exit status.

The key principle is this: when you define an oracle, the target program must be instrumented or configured to detect and report the kind of illegal behavior the oracle is designed to capture.

### 5.3.6 Fuzzing Campaign

We have introduced the core components of the fuzzing process. Now, it's time to put all the pieces together and illustrate the complete workflow.

1. The fuzzer initializes shared memory for communication between itself and the runtime.
2. The fuzzer loads the target process (which remains blocked, waiting for input).
3. The fuzzer feeds an initial seed to the target by waking up the fork-server.
4. The target program forks, and the child process executes the main logic using the provided seed.
5. The fuzzer collects the child process's exit status and coverage information, then evaluates a score for the given seed.
6. The fuzzer checks whether the seed caused a crash. If so, it minimizes the input (as described in the next chapter) and saves it to disk.
7. The fuzzer generates a new seed and returns to step 1.

This cycle is often referred to as a fuzzing campaign.

Topics such as crash minimization and symbolic execution will be explored in detail in the following chapters.

## 5.4 Implementation

As with previous components, two parts are required: instrumentation and runtime. However, at this stage, we also need to develop an additional, standalone binary—the fuzzer itself. Let's start by exploring the instrumentation component to record basic branch coverage.

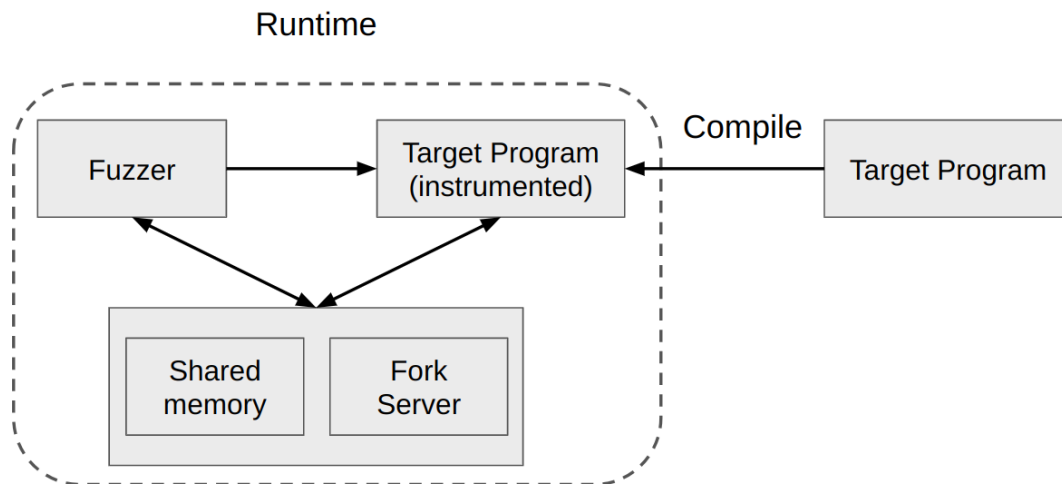


Figure 5.4: Implementation Components

### 5.4.1 Instrumentation

The instrumentation module is lightweight and primarily handles two tasks: initializing the fork-server and measuring code coverage.

#### 5.4.1.1 Fork-server Initialization

Fork-server initialization sets up the target program to receive signals from the fuzzer. Once initialized, it simply blocks and waits until a signal is received. That's all it does.

This initialization routine in the runtime is triggered via a constructor function.

```

1  // instrument/src/fuzzer.rs Rust
2  pub fn build_fuzzer_init<'ctx>(
3      context: &'ctx Context,
4      module: &Module<'ctx>,
5      builder: &Builder<'ctx>,
6  ) -> Result<FunctionValue<'ctx>> {
7      // __fuzzer_forkserver_init() - initialize fuzzer
8      let init_func_typ = context.void_type().fn_type(&[], false);
9      let init_fn = module.add_function(FUZZER_FORKSERVER_INIT,
10         init_func_typ, None);
11
12      // Create a constructor function to initialize fuzzer
13      let constructor = module.add_function(FUZZER_MODULE_INIT,
14         init_func_typ, None);
15      let entry = context.append_basic_block(constructor,
16         FUZZER_INIT_ENTRY);
17      builder.position_at_end(entry);
18      builder.build_call(init_fn, &[], "");
19      builder.build_return(None);

```



```

17     Ok(constructor)
18 }

```

It defines and creates the function `__fuzzer_module_init()`, which consists of a single basic block that calls the runtime function `__fuzzer_forkserver_init()`. The following LLVM IR snippet is generated by the `build_fuzzer_init()` function.

```

1 declare void @__fuzzer_forkserver_init()
2
3 define void @__fuzzer_module_init() {
4   __fuzzer_init_entry:
5     call void @__fuzzer_forkserver_init()
6     ret void
7 }

```

LLVM-IR

#### 5.4.1.2 Basic Block(Edge) Coverage

The remaining task is to instrument the runtime function to record coverage for each basic block.

```

1  // instrument/src/fuzzer.rs
2  fn get_rand_value<'ctx>(context: &'ctx Context) -> IntValue<'ctx> {
3      let (lower, _) = Uuid::new_v4().as_u64_pair();
4      let i64_typ = context.i64_type();
5      i64_typ.const_int(lower, false)
6  }
7
8  #[derive(Default)]
9  pub struct FuzzModule {}
10
11 impl InstrumentModule for FuzzModule {
12     fn instrument<'ctx>(
13         &self,
14         context: &'ctx Context,
15         module: &Module<'ctx>,
16         builder: &Builder<'ctx>,
17     ) -> Result<()> {
18         let constructor = build_fuzzer_init(context, module,
19             builder)?;
20         build_ctros(context, module, constructor)?;
21
22         let module_filename =
23             cstr_to_str(module.get_source_file_name());
24         let funcs: Vec<_> = module.get_functions().collect();
25         for func in funcs {

```

Rust

```

24      // Skip funcs without bodies or those we've added
25      if can_skip_instrument(&func) {
26          continue;
27      }
28      let mut instrumented_blks = HashSet::new();
29      for basic_blk in func.get_basic_blocks() {
30          if instrumented_blks.contains(&basic_blk) {
31              continue;
32          }
33          if let Some(first_instr) =
34              basic_blk.get_first_instruction() {
35              let mut instrument_pos = first_instr;
36              match first_instr.get_opcode() {
37                  InstructionOpcode::LandingPad |
38                  InstructionOpcode::Phi => {
39                      if let Some(second_instr) =
40                          first_instr.get_next_instruction() {
41                          instrument_pos = second_instr;
42                      }
43                  }
44                  _ => {}
45              }
46              if let Some(instr_filename) =
47                  get_instr_filename(&instrument_pos) {
48                  if instr_filename != module_filename {
49                      continue;
50                  }
51                  builder.position_before(&instrument_pos);
52                  build_trace_edge(context, module, builder,
53                      get_rand_value(context))?;
54              }
55              instrumented_blks.insert(basic_blk);
56          }
57      }
58
59      // Verify instrumented IRs
60      module_verify(module)
61  }
62 }
63
64 // instrument/src/inkwell_intrinsic.rs
65 fn get_trace_edge<'ctx>(context: &'ctx Context, module:
66     &Module<'ctx>) -> FunctionValue<'ctx> {

```

```

61     match get_func(module, FUZZER_TRACE_EDGE) {
62         Some(func) => func,
63         None => {
64             let trace_edge = context
65                 .void_type()
66                 .fn_type(&[context.i64_type().into()], false);
67             module.add_function(FUZZER_TRACE_EDGE, trace_edge, None)
68         }
69     }
70 }
71
72 pub fn build_fuzzer_init<'ctx>(
73     context: &'ctx Context,
74     module: &Module<'ctx>,
75     builder: &Builder<'ctx>,
76 ) -> Result<FunctionValue<'ctx>> {
77     // __fuzzer_forkserver_init() - initialize fuzzer
78     let init_func_typ = context.void_type().fn_type(&[], false);
79     let init_fn = module.add_function(FUZZER_FORKSERVER_INIT,
80                                     init_func_typ, None);
81
82     // Create a constructor function to initialize fuzzer
83     let constructor = module.add_function(FUZZER_MODULE_INIT,
84                                         init_func_typ, None);
85     let entry = context.append_basic_block(constructor,
86                                         FUZZER_INIT_ENTRY);
87     builder.position_at_end(entry);
88     builder.build_call(init_fn, &[], "");
89     builder.build_return(None);
90     Ok(constructor)
91 }

```

The implementation is straightforward. It inserts a call to the `__fuzzer_trace_edg()` function at lines 48–49, which is handled by the runtime.

The resulting instrumented LLVM IR looks as follows:

```

1  bb1:
2      ...
3      call void @__fuzzer_trace_edge(i64 @-2510630191682466935)
4      ...
5
6  bb2:
7      ...

```

LLVM-IR

```

8    call void @__fuzzer_trace_edge(i64 6646994726620055515)
9    ...
10
11 ...

```

The argument passed to `__fuzzer_trace_edge()` corresponds to `cur_location = <COMPILE_TIME_RANDOM>`. This means that for each basic block, a compile-time generated random value is assigned and used as the current location during runtime. With the instrumentation complete, we now move on to implementing the runtime library.

### 5.4.2 Runtime

There are only two external functions.

```

1  // fuzzer_runtime/src/runtime.rs Rust
2  #[no_mangle]
3  pub extern "C" fn __fuzzer_trace_edge(cur_loc: i64) {
4      EDGE_COVERAGE.lock().unwrap().trace_edge(cur_loc);
5  }
6
7  #[no_mangle]
8  pub extern "C" fn __fuzzer_forkserver_init() {
9      EDGE_COVERAGE.lock().unwrap().read_wakeup();
10 }

```

Before looking into the body of `read_wakeup()`, let's take a look the data structure.

```

1  // fuzzer_runtime/src/coverage.rs Rust
2  lazy_static::lazy_static! {
3      pub static ref EDGE_COVERAGE: Arc<Mutex<EdgeCoverage>> =
4          Arc::new(Mutex::new(EdgeCoverage::new()));
5  }
6
7  pub struct EdgeCoverage {
8      shm: Option<MmapMut>,
9      shm_size: Option<usize>,
10     aux: Option<MmapMut>,
11     fork_server_host: Option<i32>,
12     fork_server_runtime: Option<i32>,
13 }
14
15 impl Default for EdgeCoverage {
16     fn default() -> Self {
17         Self::new()
18     }
19 }

```

```

18 }
19
20 impl EdgeCoverage {
21     pub fn new() -> Self {
22         let (shm_mmap, shm_size) = init_shm("SHM_ID", "SHM_SIZE");
23         let (shm_aux_mmap, _) = init_shm("SHM_AUX_ID",
24             "SHM_AUX_SIZE");
25         Self {
26             shm: shm_mmap,
27             shm_size,
28             aux: shm_aux_mmap,
29             fork_server_host:
30                 init_forkserver_fd("FORK_SERVER_HOST"),
31             fork_server_runtime:
32                 init_forkserver_fd("FORK_SERVER_RUNTIME"),
33         }
34     }
35
36     /// only used for test
37     pub fn init(&mut self) {
38         let new_edge_cov = Self::new();
39         self.shm = new_edge_cov.shm;
40         self.shm_size = new_edge_cov.shm_size;
41         self.aux = new_edge_cov.aux;
42         self.fork_server_host = new_edge_cov.fork_server_host;
43         self.fork_server_runtime = new_edge_cov.fork_server_runtime;
44     }
45
46     ...
47 }

```

EdgeCoverage is implemented as a singleton object. The `shm` and `aux` variables are file descriptors pointing to shared memory regions. These shared memory segments are mapped to the `/tmp` directory using `mmap()` [23]. The actual creation and mapping of this virtual memory space are delegated to the fuzzer.

The `shm` region stores execution data—specifically, which basic blocks have been executed and how many times. The `aux` region holds auxiliary metadata, such as the total number of edges, which edges are discovered and what’s new finding.

The `fork_server_host` and `fork_server_runtime` are file descriptors used for event notification between the fuzzer and the runtime. Signaling between the two is performed via these event file descriptors, using mechanisms such as `eventfd` [24].

```

1 // fuzz_runtime/src/internal.rs
2 pub fn set_mmap(id: &str) -> Option<MmapMut> {

```



```

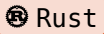
3      match env::var(id) {
4          Ok(path) => {
5              let file = OpenOptions::new()
6                  .read(true)
7                  .write(true)
8                  .open(path)
9                  .expect("Failed to open SHM file");
10             let mmap = unsafe { MmapMut::map_mut(&file).expect("mmap
11                 failed in runtime") };
12             Some(mmap)
13         }
14         Err(_) => None,
15     }
16
17 pub fn init_shm(id: &str, size_id: &str) -> (Option<MmapMut>,
18     Option<usize>) {
19     let ret = set_mmap(id);
20     let mut shm_size = None;
21     if let Ok(size) = env::var(size_id) {
22         shm_size = Some(size.parse:::<usize>().unwrap());
23     }
24     (ret, shm_size)
25 }
26
27 pub fn init_forkserver_fd(id: &str) -> Option<i32> {
28     if let Ok(fd) = env::var(id) {
29         return Some(fd.parse:::<i32>().unwrap());
30     }
31     None
32 }

```

Initializing these file descriptors is straightforward.

As instrumented, the `__fuzzer_forkserver_init()` function sets up a loop that waits for a signal from the fuzzer to begin execution.

```

1  // fuzz_runtime/src/coverage.rs 
2  pub const PROCESS_EXIT_NORMAL: u64 = 1;
3  pub fn read_wakeup(&self) {
4      if let Some(fd) = self.fork_server_runtime {
5          loop {
6              let mut host_sent: u64 = 0;
7              unsafe {

```

```
8         read(
9             fd,
10             &mut host_sent as *mut _ as *mut c_void,
11             mem::size_of::<u64>(),
12         );
13     }
14     if host_sent == 9999999999 {
15         // exit signal
16         unsafe {
17             libc::exit(0);
18         }
19     }
20     let pid = unsafe { libc::fork() };
21     if pid == 0 {
22         // child process
23         return; // run target program's main logic
24     } else {
25         // parent process
26         let (tx, rx) = mpsc::channel();
27         thread::spawn(
28             move || match
29                 rx.recv_timeout(Duration::from_secs(host_sent))
30             {
31                 Ok(_) => {}
32                 Err(mpsc::RecvTimeoutError::Timeout) =>
33                     unsafe {
34                         kill(pid, SIGKILL);
35                     }
36                 _ => {}
37             },
38         );
39         let mut status: i32 = 0;
40         unsafe {
41             waitpid(pid, &mut status, 0);
42             tx.send(()).ok();
43         }
44         status = if status != SIGKILL {
45             WEXITSTATUS(status)
46         } else {
47             status
48         };
49     }
```

```

46         self.notify_process_exit(status.try_into().unwrap());
47     }
48 }
49 }
50 }

```

The `read_wakeup()` function enters an infinite loop, performing a blocking read to wait for a wakeup signal from the fuzzer. If the fuzzer sends the value 9999999999, it indicates that the fuzzing process is complete, and the runtime can safely release resources. Otherwise, the received value is treated as a timeout, which sets a time limit for the target program's execution.

Once the signal is received, the runtime forks the target process using `libc::fork()`. The child process executes the target program's main logic using the provided fuzzing input (which fed from `stdin` or read from file), while the parent process waits for the child to finish. If the child process does not terminate within the specified timeout, it is forcibly killed.

```

1  // fuzz_runtime/src/coverage.rs Rust
2  fn notify_process_exit(&self, status: u64) {
3      if let Some(fd) = self.fork_server_host {
4          // do not send zero value
5          let status = if status == 0 {
6              PROCESS_EXIT_NORMAL
7          } else {
8              status
9          };
10         let bytes = status.to_ne_bytes();
11         unsafe { write(fd, bytes.as_ptr() as *const _,
12             bytes.len()) };
13     }
14 }

```

```

1  // fuzz_runtime/src/internal.rs Rust
2  pub fn read_u64(mem: &[u8], start: usize) -> u64 {
3      let mut slice = [0u8; 8];
4      slice.copy_from_slice(&mem[start..start + 8]);
5      u64::from_le_bytes(slice)
6  }
7
8  pub fn read_u128(mem: &[u8], start: usize) -> u128 {
9      let mut slice = [0u8; 16];
10     slice.copy_from_slice(&mem[start..start + 16]);
11     u128::from_le_bytes(slice)

```



```

12 }
13
14 pub fn read_cov_report(mem: &[u8]) -> Vec<CovReport> {
15     let size = read_u64(mem, 0) as usize;
16     let mut slice = vec![0u8; size];
17     slice.copy_from_slice(&mem[8..8 + size]);
18     if let Ok(des) = bincode::deserialize(&slice) {
19         des
20     } else {
21         vec![]
22     }
23 }
24
25 pub fn write_u64(mem: &mut [u8], start: usize, v: usize) {
26     mem[start..start + 8].copy_from_slice(&v.to_le_bytes())
27 }
28
29 pub fn write_u128(mem: &mut [u8], start: usize, v: u128) {
30     mem[start..start + 16].copy_from_slice(&v.to_le_bytes())
31 }

```

Finally, the parent process checks the exit status of the child process and sends it back to the fuzzer. This status indicates whether the program terminated abnormally, allowing the fuzzer to determine whether a crash has occurred.

The second function is coverage measurement, which is handled by `trace_edge()`.

```

1 // fuzzer_runtime/src/coverage.rs Rust
2
3 const PREV_LOC_IDX: usize = 0; // 16 byte [0..15]
4 pub const NEW_COVERAGES: usize = PREV_LOC_IDX + 16; // 8 byte
   [16..23]
5 pub const VISIT_EDGES: usize = NEW_COVERAGES + 8; // 8 bytes
   [24..31]
6 pub const VISIT_MARK: usize = VISIT_EDGES + 8; // 8 bytes [32..39]
7 pub const VISIT_EDGES_INDICIES: usize = VISIT_MARK + 8; // n bytes
   [40..]
8 pub const VISIT_EDGES_INDEX_SIZE: usize = 8;
9
10 pub fn trace_edge(&mut self, cur_loc: i64) {
11     if let (Some(ref mut shm), Some(shm_size), Some(ref mut
   shm_aux)) =
12         (&mut self.shm, self.shm_size, &mut self.aux)
13     {
14         let cur_loc = cur_loc as u128;

```

```

15     let prev_loc = read_u128(shm_aux, PREV_LOC_IDX);
16     let edge = (cur_loc ^ prev_loc) as usize % shm_size;
17     write_u128(shm_aux, PREV_LOC_IDX, cur_loc >> 1);
18
19     let visited_cnt = read_u64(shm_aux, VISIT_MARK) as usize;
20     let already_visited = (0..visited_cnt).any(|i| {
21         let stored_edge =
22             read_u64(shm_aux, VISIT_EDGES_INDICIES + (i *
23                 VISIT_EDGES_INDEX_SIZE));
24         stored_edge as usize == edge
25     });
26     if !already_visited {
27         write_u64(
28             shm_aux,
29             VISIT_EDGES_INDICIES + (visited_cnt *
30                 VISIT_EDGES_INDEX_SIZE),
31             edge,
32         );
33         write_u64(shm_aux, VISIT_MARK, visited_cnt + 1);
34     }
35     if shm[edge] == 0 {
36         write_u64(shm_aux, NEW_COVERAGES, 1);
37         let edges = read_u64(shm_aux, VISIT_EDGES) as usize;
38         write_u64(shm_aux, VISIT_EDGES, edges + 1);
39     }
40     shm[edge] = shm[edge].saturating_add(1);
41 }

```

We instrument a call to the `__fuzzer_trace_edge()` runtime function at every basic block. This function, in turn, calls `trace_edge` to process each executed edge.

The first step within `trace_edge()` is to calculate a unique edge index. To do this, we use `shm_aux` to store the previous basic block's location. This allows us to represent the transition (edge) between the previous and current block.

Next, we check if this specific edge has been visited before during the current execution. If it's a new edge for this run, we mark it as visited and increment a counter for the total number of unique edges visited. This information is crucial for the fuzzer, as it needs to know which edges were covered by the current seed and how frequently. This helps in seed evaluation through bucketing, a process where the fuzzer assesses the novelty and impact of each input. To accurately track per-run coverage, we record all visited edges for each target program invocation and clear this specific set of visited edges once the program exits. This ensures that for every new run, the fuzzer has a precise understanding of the edges covered by that specific input.

Finally, the overall visit counts for all edges are managed within `shm` (shared memory). It's critical to note that `shm` is not cleared after each program execution. This persistence

allows the fuzzer to maintain a cumulative record of how frequently each edge has been hit across all fuzzing runs. If `shm` were cleared, the fuzzer would treat every edge as “new” in every single execution.

For example, imagine the fuzzer discovers a new code path when it processes Input A. Now, in the next run, the fuzzer mutates Input A to create Input B and feeds it to the program. If Input B happens to traverse the exact same code path as Input A, there’s no “new” discovery in terms of unique paths for this particular run. **This is precisely why the edge count must persist across runs.** If it didn’t, the fuzzer would repeatedly register already-found paths as “new,” significantly hindering its ability to efficiently explore truly novel code.

### 5.4.3 Fuzzer

We’ve successfully implemented both the instrumentation module and its corresponding runtime handler. Our next and final step is to implement the fuzzer itself. Let’s begin by examining its `main()` function.

#### 5.4.3.1 Initialization

```

1  // fuzzer/src/mmap.rs Rust
2  pub const SHM_PATH: &str = "/tmp/fuzzer_shared_mem";
3  pub const SHM_AUX_PATH: &str = "/tmp/fuzzer_shared_aux_mem";
4  pub const SHM_COV_PATH: &str = "/tmp/fuzzer_shared_cov_mem";
5  pub const SHM_SIZE: usize = 1 << 16;
6  pub const SHM_AUX_SIZE: usize = 1 << 20;
7  pub const SHM_COV_SIZE: usize = 1 << 13;
8
9  // fuzzer/src/main.rs
10 fn main() -> Result<()> {
11     let (pgm_path, seed_dir_path, input_typ) = get_args()?;
12     let init_seeds = SeedPool::new(&seed_dir_path);
13     let (tx, rx) = mpsc::channel();
14     let mut fuzzer = Fuzzer::new(
15         SHM_PATH,
16         SHM_SIZE,
17         SHM_AUX_PATH,
18         SHM_AUX_SIZE,
19         SHM_COV_PATH,
20         SHM_COV_SIZE,
21         init_seeds,
22         input_typ,
23         tx,
24     );
25     let handle = std::thread::spawn(move || {
26         let _ = run_ui(rx);

```

```

27     });
28     let _ = fuzzer.run(&pgm_path);
29     handle.join().unwrap();
30     Ok(())
31 }

```

The `main()` function initializes:

- CLI: Takes program arguments
- SeedPool: Construct a seed pool object which is detailed later
- Fuzzer: Construct a fuzzer object
- UI: Configure fuzzer UI, which is detailed later
- Fuzzing campaign: Run fuzzing process

We won't include the command-line interface (CLI) implementation here; you can find it in the source code repository. The fuzzer's initialization primarily involves setting up shared memory.

```

1  // fuzzer/src/fuzzer.rs Rust
2  pub struct Fuzzer {
3      shm: SHM,
4      shm_aux: SHM,
5      shm_cov: SHM,
6      pub forkserver_host: i32,
7      pub forkserver_runtime: i32,
8      seeds: SeedPool,
9      new_paths: usize,
10     input_typ: FuzzInput,
11     seed_file_path: String, // if fuzz input is `ProgramArgument`,
                             // we write seed into this file path
12     crashes: HashSet<Seed>,
13     tx: mpsc::Sender<FuzzShot>,
14     fuzz_terminate: bool,
15 }
16
17 impl Fuzzer {
18     pub fn new(
19         shm_path: &str,
20         shm_size: usize,
21         shm_aux_path: &str,
22         shm_aux_size: usize,
23         shm_cov_path: &str,
24         shm_cov_size: usize,
25         init_seeds: SeedPool,
26         input_typ: FuzzInput,

```

```

27         tx: mpsc::Sender<FuzzShot>,
28     ) -> Self {
29         let shm = SHM::new(shm_path, shm_size);
30         let shm_aux = SHM::new(shm_aux_path, shm_aux_size);
31         let shm_cov = SHM::new(shm_cov_path, shm_cov_size);
32         let host_efd: RawFd = unsafe { eventfd(0, 0) };
33         let runtime_efd: RawFd = unsafe { eventfd(0, 0) };
34         Self {
35             shm,
36             shm_aux,
37             shm_cov,
38             forkserver_host: host_efd,
39             forkserver_runtime: runtime_efd,
40             seeds: init_seeds,
41             new_paths: 0,
42             input_typ,
43             seed_file_path: format!("{}", Uuid::new_v4()),
44             crashes: HashSet::new(),
45             tx,
46             fuzz_terminate: false,
47         }
48     }
49 }

```

Shared memory struct is defined as:

```

1  // fuzzer/src/mmap.rs Rust
2  pub struct SHM {
3      mem: MmapMut,
4      mem_path: String,
5      mem_size: usize,
6  }
7
8  impl SHM {
9      pub fn new(shm_path: &str, shm_size: usize) -> Self {
10         // 1. Create shared memory file
11         let _ = fs::remove_file(shm_path); // clean up any previous
            file
12         let file = OpenOptions::new()
13             .read(true)
14             .write(true)
15             .create(true)

```

```

16         .open(shm_path)
17         .unwrap();
18     file.set_len(shm_size as u64).unwrap();
19
20     // 2. Memory-map the file
21     let m = unsafe { MmapMut::map_mut(&file).unwrap() };
22     let mut mmap = Self {
23         mem: m,
24         mem_path: shm_path.to_string(),
25         mem_size: shm_size,
26     };
27
28     // 3. Zeroize
29     mmap.zeroize();
30     mmap
31 }
32
33 pub fn mut_mem(&mut self) -> &mut MmapMut {
34     &mut self.mem
35 }
36
37 pub fn mem(&self) -> &MmapMut {
38     &self.mem
39 }
40
41 pub fn path(&self) -> &str {
42     &self.mem_path
43 }
44
45 pub fn size(&self) -> usize {
46     self.mem_size
47 }
48
49 pub fn zeroize(&mut self) {
50     self.mem[..].fill(0);
51 }
52 }

```

#### 5.4.3.2 Seed Initialization

There are two structures: Seed and SeedPool.

The struct Seed contains two important data:

```

1  #[derive(Debug, Clone)]
2  pub struct Seed {
3      input: Vec<u8>,
4      score: u64,
5  }

```

 Rust

The input represents the seed data, and score is its evaluated value, determined by the fuzzer.

Seeds are managed by the `SeedPool` struct. This struct internally utilizes a `BTreeSet`, which means any type used within it must implement the `Hash` trait. Additionally, the `mutate()` function is available to perform modifications on a given seed.

The `SeedPool` offers two primary methods: `add_seed()` to incorporate new seeds, and `pop_seed()` to retrieve the most promising seed from the pool.

```

1  // fuzzer/src/seed.rs
2  const MAX_SEED_LEN: usize = 1000;
3  const CRASH_OUTPUT_DIR: &str = "crashes";
4
5  #[derive(Debug, Clone)]
6  pub struct Seed {
7      input: Vec<u8>,
8      score: u64,
9  }
10
11  impl Seed {
12      pub fn new(input: Vec<u8>, score: u64) -> Self {
13          Self { input, score }
14      }
15
16      pub fn mutate(&mut self) -> Vec<MutateResult> {
17          mutator::mutate(&mut self.input)
18      }
19
20      pub fn get_input(&self) -> &[u8] {
21          &self.input
22      }
23
24      pub fn get_score(&self) -> u64 {
25          self.score
26      }
27
28      pub fn to_hex(&self) -> String {
29          self.get_input()

```

 Rust

```
30         .iter()
31         .map(|b| format!("{:02X}", b))
32         .collect::<Vec<_>>()
33         .join("")
34     }
35
36     pub fn to_file(&self, cnt: usize) {
37         fs::create_dir_all(CRASH_OUTPUT_DIR).unwrap();
38         let mut file_path = PathBuf::from(CRASH_OUTPUT_DIR);
39         file_path.push(format!("{}.crash", cnt));
40         fs::write(&file_path, self.input.clone()).unwrap();
41     }
42
43     pub fn set_score(&mut self, v: u64) {
44         self.score = v;
45     }
46 }
47
48 impl Ord for Seed {
49     fn cmp(&self, other: &Self) -> std::cmp::Ordering {
50         self.score
51             .cmp(&other.score)
52             .then_with(|| self.input.cmp(&other.input))
53     }
54 }
55
56 impl PartialOrd for Seed {
57     fn partial_cmp(&self, other: &Self) ->
58         Option<std::cmp::Ordering> {
59         Some(self.cmp(other))
60     }
61
62 impl PartialEq for Seed {
63     fn eq(&self, other: &Self) -> bool {
64         self.input == other.input && self.score == other.score
65     }
66 }
67
68 impl Eq for Seed {}
69
```



```
70 impl Hash for Seed {
71     fn hash<H: Hasher>(&self, state: &mut H) {
72         self.input.hash(state);
73         self.score.hash(state);
74     }
75 }
76
77 pub struct SeedPool {
78     seeds: BTreeSet<Seed>,
79 }
80
81 impl SeedPool {
82     pub fn new(seed_dir: &str) -> Self {
83         let mut btree = BTreeSet::new();
84         let init_seeds = read_seed_dir(seed_dir).unwrap();
85         for init_seed in init_seeds {
86             btree.insert(Seed::new(init_seed, 0));
87         }
88         Self { seeds: btree }
89     }
90
91     pub fn add_seed(&mut self, seed: Seed) {
92         if self.seeds.len() > MAX_SEED_LEN {
93             if let Some(min_seed) = self.get_min_score_seed() {
94                 self.seeds.remove(&min_seed);
95             }
96         }
97         self.seeds.insert(seed);
98     }
99
100     fn get_min_score_seed(&self) -> Option<Seed> {
101         self.seeds.iter().next().cloned()
102     }
103
104     fn get_max_score_seed(&self) -> Option<Seed> {
105         self.seeds.iter().next_back().cloned()
106     }
107
108     pub fn pop_seed(&mut self) -> Option<Seed> {
109         if let Some(seed) = self.get_max_score_seed() {
110             self.seeds.remove(&seed);
```

```

111         Some(seed)
112     } else {
113         None
114     }
115 }
116
117 pub fn is_empty(&self) -> bool {
118     self.seeds.len() == 0
119 }
120
121 pub fn len(&self) -> usize {
122     self.seeds.len()
123 }
124 }

```

### 5.4.3.3 Target Program Initialization

Initially, the target program is set up using the Command crate. During this setup, several environment variables are notably configured.

- LD\_LIBRARY\_PATH: Sets current directory as library path, meaning that fuzzer runtime should locate within current directory
- SHM\_XXX: Shared memory-related values such as file path
- FORK\_SERVER\_XXX: Event notification related values

It's worth noting that stdout and stderr are piped. This allows us to intercept the target program's terminal output and render these outputs directly within our fuzzer UI.

```

1  // fuzzer/src/fuzzer.rs Rust
2  pub fn run(&mut self, program_file: &str) -> Result<FuzzResult> {
3      let mut cmd = Command::new(program_file);
4      let child_process_cmd = cmd
5          .env("LD_LIBRARY_PATH", ".")
6          .env("SHM_ID", self.shm.path())
7          .env("SHM_AUX_ID", self.shm_aux.path())
8          .env("SHM_SIZE", format!("{}", &self.shm.size()))
9          .env("SHM_AUX_SIZE", format!("{}", &self.shm_aux.size()))
10         .env("SHM_COV_ID", self.shm_cov.path())
11         .env("SHM_COV_SIZE", format!("{}", &self.shm_cov.size()))
12         .env("FORK_SERVER_HOST", format!("{}",
13             self.forkserver_host))
14         .env(
15             "FORK_SERVER_RUNTIME",
16             format!("{}", self.forkserver_runtime),
17         )
18         .stdout(Stdio::piped())

```

```

18         .stderr(Stdio::piped());
19
20     if self.input_typ == FuzzInput::ProgramArgument {
21         child_process_cmd.args(vec! [&self.seed_file_path]);
22     } else {
23         child_process_cmd.stdin(Stdio::piped());
24     }
25
26     let mut child_process = child_process_cmd.spawn()?;
27     let mut child_stdin = if self.input_typ == FuzzInput::Stdin {
28         Some(child_process.stdin.take().unwrap())
29     } else {
30         None
31     };
32
33     ...
34     ...
35
36 }

```

At line 26, we load the target program, which will then enter a blocked state.

#### 5.4.3.4 Harness

As previously explained, the target program needs to implement a harness to interact with the fuzzer. In this book, we'll use two methods for this interaction: `stdin` and `file`.

If `stdin` is chosen, we intercept the input at lines 27-29. The `feed_seed()` function then determines which harness is in use. If it's `stdin`, the function writes a seed directly to the `stdin` file descriptor; otherwise, it writes the seed into a file.

```

1  fn feed_seed(&self, child_stdin: &mut Option<ChildStdin>,
2  seed: &Seed) -> Result<()> {
3      if let Some(ref mut stdin) = child_stdin {
4          stdin.write_all(seed.get_input())?;
5      }
6      if self.input_typ == FuzzInput::ProgramArgument {
7          write_seed(&self.seed_file_path, seed.get_input())?;
8      }
9      Ok(())
10 }

```

#### 5.4.3.5 Fuzzing Campaign

The entire fuzzing process is attached below:

```

1  pub fn run(&mut self, program_file: &str) ->
2  Result<FuzzResult> {

```

```
2     let mut cmd = Command::new(program_file);
3     let child_process_cmd = cmd
4         .env("LD_LIBRARY_PATH", ".")
5         .env("SHM_ID", self.shm.path())
6         .env("SHM_AUX_ID", self.shm_aux.path())
7         .env("SHM_SIZE", format!("{}", &self.shm.size()))
8         .env("SHM_AUX_SIZE", format!("{}", &self.shm_aux.size()))
9         .env("SHM_COV_ID", self.shm_cov.path())
10        .env("SHM_COV_SIZE", format!("{}", &self.shm_cov.size()))
11        .env("FORK_SERVER_HOST", format!("{}",
12            self.forkserver_host))
13        .env(
14            "FORK_SERVER_RUNTIME",
15            format!("{}", self.forkserver_runtime),
16        )
17        .stdout(Stdio::piped())
18        .stderr(Stdio::piped());
19
20    if self.input_typ == FuzzInput::ProgramArgument {
21        child_process_cmd.args(vec![&self.seed_file_path]);
22    } else {
23        child_process_cmd.stdin(Stdio::piped());
24    }
25
26    let mut child_process = child_process_cmd.spawn()?;
27    let mut child_stdin = if self.input_typ == FuzzInput::Stdin {
28        Some(child_process.stdin.take().unwrap())
29    } else {
30        None
31    };
32
33    self.pgm_output_reader(
34        child_process.stdout.take().unwrap(),
35        child_process.stderr.take().unwrap(),
36    );
37
38    let mut init_set_timeout = false;
39    let mut timeout = Duration::new(9999, 0);
40    let fuzzer_started = Instant::now();
41
42    let mut loop_cnt = 0;
43    if self.is_seed_empty() {
```

```
42         return Ok(FuzzResult::AllSeedConsumed(loop_cnt));
43     }
44     let mut seed = self.pop_seed().unwrap();
45
46     loop {
47         if self.fuzz_terminate {
48             self.terminate();
49             self.send(FuzzShot::Terminated);
50             return Ok(FuzzResult::UserTerminated);
51         }
52         loop_cnt += 1;
53         self.feed_seed(&mut child_stdin, &seed)?;
54
55         self.wakeup_forkserver(HostSend::Wakeup(timeout.as_secs()));
56         let target_started = Instant::now();
57         let status = self.wait_forkserver();
58
59         let elapsed = target_started.elapsed();
60         if !init_set_timeout {
61             timeout = elapsed * INITIAL_TIMEOUT_UPPER_BOUND;
62             init_set_timeout = true;
63         }
64         // evaluate seed and add it
65         let (visit_edges, score) = self.eval_seed(status);
66         self.debug(
67             loop_cnt,
68             &seed,
69             score,
70             timeout,
71             elapsed,
72             fuzzer_started.elapsed(),
73         );
74         self.clear_new_coverage();
75         self.clear_visited_edges();
76         if self.is_seed_empty() {
77             self.terminate();
78             self.send(FuzzShot::Terminated);
79             return Ok(FuzzResult::AllSeedConsumed(loop_cnt));
80         }
81         // crash found
```

```

82         if self.is_crash(status) {
83             if let Ok(minimized) = self.ddmin(&mut child_stdin,
                timeout, &seed) {
84                 self.crashes.insert(minimized.clone());
85                 minimized.to_file(self.crashes.len());
86                 self.send(FuzzShot::Crash(CrashInfo::new(
87                     self.crashes.len(),
88                     seed,
89                     minimized,
90                     )));
91             }
92         } else if score > 0 {
93             // re-evaluate score after execution
94             seed.set_score(score);
95             self.add_seed(seed.clone());
96             let cur_seed = seed.clone();
97             seed.mutate();
98             seed.set_score(score.saturating_add(1));
99             self.add_seed(seed.clone());
100
101             self.send(FuzzShot::SeedInfo(FuzzerSeed::new(
102                 self.seeds.len(),
103                 cur_seed,
104                 seed,
105                 visit_edges,
106                 self.new_paths,
107                 )));
108         }
109         seed = self.pop_seed().unwrap();
110     }
111 }

```

At lines 36-38, we initialize the timeout and record the fuzzer's start time. One seed is popped from the seed pool between lines 41 and 44. The main fuzzing loop begins at line 46. If the fuzzer should terminate, it exits the function at lines 47-51. This termination condition becomes true when the UI is exited (e.g., by typing q), a signal handled by the UI implementation.

We feed a seed input at line 53 and then send a wakeup signal to the fuzzer runtime at line 54. Following this, the target program will fork and execute its main logic. The fuzzing input will be delivered via the harness.

After the target program begins execution, the fuzzer waits for its completion. Once the target program finishes, it sends its exit status back to the fuzzer.

```

2  pub fn wait_forkserver(&mut self) -> i32 {
3      let mut status: u64 = 0;
4      unsafe {
5          read(
6              self.forkserver_host,
7              &mut status as *mut _ as *mut c_void,
8              mem::size_of::<u64>(),
9          );
10     }
11     status as i32
12 }

```

At line 64, the given seed is evaluated, producing a score value. If new coverage is discovered, the seed receives a higher score; otherwise, its score is determined through a process called bucketing.

```

1  // fuzzer/src/fuzzer.rs 🔒 Rust
2  pub fn eval_seed(&mut self, status: i32) -> (u64, u64) {
3      if status == SIGKILL {
4          return (0, 0); // give zero score for hangs
5      }
6      let visit_edges = read_u64(self.shm_aux.mem(), VISIT_MARK);
7      // if new coverage found, give max score, otherwise give higher
8      // score if the path is rare
9      if self.is_new_coverage() {
10         self.new_paths += 1;
11         return (visit_edges, u64::MAX);
12     }
13     let mut score = 0;
14     for i in 0..visit_edges {
15         let edge = read_u64(
16             self.shm_aux.mem(),
17             VISIT_EDGES_INDICIES + (i as usize) *
18             VISIT_EDGES_INDEX_SIZE,
19         ) as usize;
20         if edge != 0 {
21             score += get_score(self.shm.mem()[edge]) as u64;
22         }
23     }
24     (visit_edges, score)
25 }

```

Bucketing is implemented as follows:

```

1 // fuzzer/src/bucket.rs Rust
2 const BUCKET_MAX_VALUE: u8 = 8;
3 const COUNT_CLASS_LOOKUP: [u8; 256] = {
4     let mut table = [0u8; 256];
5     let mut i = 0;
6     while i < 256 {
7         table[i] = match i {
8             0..=1 => 1,
9             2 => 2,
10            3 => 3,
11            4..=7 => 4,
12            8..=15 => 5,
13            16..=31 => 6,
14            32..=127 => 7,
15            _ => BUCKET_MAX_VALUE,
16        };
17        i += 1;
18    }
19    table
20 };
21
22 /// lower value is more valuable because it has been not rarely
23   exercised
24 fn to_bucket(value: u8) -> u8 {
25     COUNT_CLASS_LOOKUP[value as usize]
26 }
27 pub fn get_score(value: u8) -> u8 {
28     let score = to_bucket(value);
29     (BUCKET_MAX_VALUE + 1) - score
30 }

```

As described earlier, a lower hit count is more valuable, while a higher hit count will be assigned a lower score.

Once the score evaluation is complete, the coverage footprint — specifically, the number of visited edges, which edges are new, and which edges have been visited — must be cleared. This information can be reset via the `shm_aux` shared memory.

```

1 // fuzzer/src/fuzzer.rs Rust
2 pub fn clear_new_coverage(&mut self) {
3     write_u64(self.shm_aux.mut_mem(), NEW_COVERAGES, 0);
4 }
5

```



```

6  pub fn clear_visited_edges(&mut self) {
7      let visit_edges = read_u64(self.shm_aux.mem(), VISIT_MARK);
8      for i in 0..visit_edges {
9          write_u64(
10             self.shm_aux.mut_mem(),
11             VISIT_EDGES_INDICIES + (i as usize) *
12             VISIT_EDGES_INDEX_SIZE,
13             0,
14         );
15     }
16     write_u64(self.shm_aux.mut_mem(), VISIT_MARK, 0);
17 }

```

If the seed pool becomes empty, the fuzzer halts, as shown at lines 75-79. At this point, it also sends a terminate event to the runtime.

Should the current seed lead to a crash, the fuzzer will minimize the input and store it to disk (lines 82-91).

Finally, we reset the evaluated score for the current seed and add it back into the pool. We also create a copy of this seed, mutate it, and then add the mutated version to the pool. Before the loop restarts, a promising seed is selected for the next execution (line 92 - 109).

#### 5.4.3.6 Mutation

Up to this point, we've covered the fuzzer's entire workflow, but we haven't yet discussed how new seeds are generated. Our approach is to create new seeds by mutating the existing seeds. In this book, we will focus on five widely used and effective mutation methods.

- Insertion: Inserts one random byte at a random position within the seed
- Deletion: Removes a single random byte from the seed
- Change: Replaces one random byte in the seed with a new random byte
- Flip: Flips the bits of a single random byte (e.g., `seed[idx] ^= 0xFF`)
- Arithmetic: Adds a randomly generated integer (within the range of -35 to +35, as used by AFL) to a random position in the seed

The five methods were implemented as follows:

```

1  // fuzzer/src/mutator.rs Rust
2
3  use rand::Rng;
4
5  // 1. random insert
6  // 2. random delete
7  // 3. random change
8  // 4. random flip
9  // 5. arithmetic mutation
10

```

```
11  const KEEP_SEED_MIN_LEN: usize = 1;
12
13  #[derive(PartialEq, Eq, Debug)]
14  pub enum MutateResult {
15      EmptyInput,
16      SeedTooShortToDelete,
17      SeedTooShortToArithmeticMutate,
18      Done,
19  }
20
21  fn gen_new_byte() -> u8 {
22      rand::random::<u8>()
23  }
24
25  fn gen_random_idx(len: usize) -> usize {
26      rand::random::<usize>() % len
27  }
28
29  type MutatorFn = fn(&mut Vec<u8>) -> MutateResult;
30
31  const MUTATORS: &[MutatorFn] = &[insert, change, delete, flip,
    arithmetic];
32
33  fn insert(seed: &mut Vec<u8>) -> MutateResult {
34      if seed.is_empty() {
35          return MutateResult::EmptyInput;
36      }
37      let idx = gen_random_idx(seed.len());
38      seed.insert(idx, gen_new_byte());
39      MutateResult::Done
40  }
41
42  fn change(seed: &mut Vec<u8>) -> MutateResult {
43      if seed.is_empty() {
44          return MutateResult::EmptyInput;
45      }
46      let idx = gen_random_idx(seed.len());
47      seed[idx] = gen_new_byte();
48      MutateResult::Done
49  }
50
```

```
51 fn delete(seed: &mut Vec<u8>) -> MutateResult {
52     if seed.is_empty() {
53         return MutateResult::EmptyInput;
54     }
55     let idx = gen_random_idx(seed.len());
56     if seed.len() > KEEP_SEED_MIN_LEN {
57         seed.remove(idx);
58         MutateResult::Done
59     } else {
60         MutateResult::SeedTooShortToDelete
61     }
62 }
63
64 fn flip(seed: &mut Vec<u8>) -> MutateResult {
65     if seed.is_empty() {
66         return MutateResult::EmptyInput;
67     }
68     let idx = gen_random_idx(seed.len());
69     seed[idx] ^= 0xFF;
70     MutateResult::Done
71 }
72
73 fn arithmetic(seed: &mut Vec<u8>) -> MutateResult {
74     if seed.is_empty() {
75         return MutateResult::EmptyInput;
76     }
77     let mut rng = rand::thread_rng();
78     let mutate_widths = [1, 2, 4];
79     let width = match seed.len() {
80         0..=1 => return MutateResult::SeedTooShortToArithmeticMutate,
81         2..=3 => 1,
82         4 => mutate_widths[rng.gen_range(0..2)],
83         _ => mutate_widths[rng.gen_range(0..mutate_widths.len())],
84     };
85     let idx = rng.gen_range(0..=seed.len() - width);
86     // choose a random value from -35 ~ +35
87     let delta = rng.gen_range(-35i64..=35);
88     match width {
89         1 => {
90             let val = seed[idx];
```

```

91         let new_val = val.wrapping_add(delta as u8);
92         seed[idx] = new_val;
93     }
94     2 => {
95         let val = u16::from_le_bytes([seed[idx], seed[idx + 1]]);
96         let new_val = val.wrapping_add(delta as u16);
97         let bytes = new_val.to_le_bytes();
98         seed[idx..idx + 2].copy_from_slice(&bytes);
99     }
100    4 => {
101        let val = u32::from_le_bytes([seed[idx], seed[idx + 1],
102        seed[idx + 2], seed[idx + 3]]);
102        let new_val = val.wrapping_add(delta as u32);
103        let bytes = new_val.to_le_bytes();
104        seed[idx..idx + 4].copy_from_slice(&bytes);
105    }
106    _ => unreachable!(),
107 }
108 MutateResult::Done
109 }
110
111 pub fn mutate(seed: &mut Vec<u8>) -> Vec<MutateResult> {
112     // mutate 1 ~ 10% for a given seed
113     let mut rng = rand::thread_rng();
114     let max = ((seed.len() as f32) * 0.2).ceil() as usize;
115     let max = max.max(1).min(seed.len());
116     let n = rng.gen_range(1..=max);
117
118     let mut results = vec![];
119     for _ in 0..=n {
120         let idx = gen_random_idx(MUTATORS.len());
121         let result = MUTATORS[idx](seed);
122         results.push(result);
123     }
124     results
125 }

```

Each mutator is chosen randomly, and the degree of mutation has been predetermined for this book. We apply mutations ranging from 1% to 10% per seed, meaning each seed will be altered by at most 10% of its original content.

### 5.4.3.7 UI

From a functional standpoint, everything is now implemented. You'll soon see the fuzzer's user interface, which is our last piece of the puzzle: wrapping it all with an intuitive and awesome UI. Our intended UI layers are as follows:

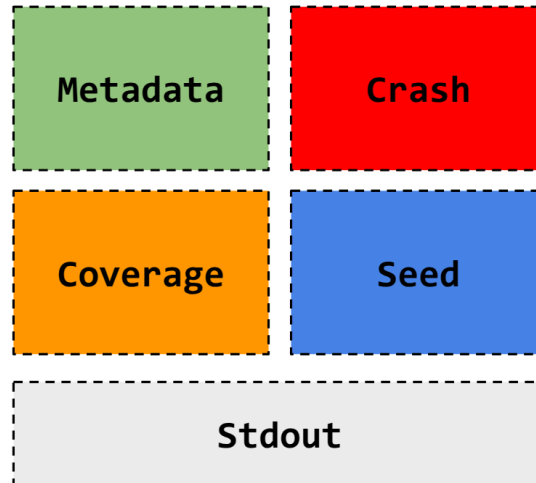


Figure 5.5: Implementation Components

There are five compartments.

- Metadata
  - Fuzzing loop counter: Increased if target program is ran once
  - Fuzzing input type: `<stdin | file>`
  - Fuzzing Timeout: Time constraint of target program invocation
  - Elapsed time of single run: An elapsed time of target program invocation
  - Uptime: Fuzzer's uptime
- Crash
  - Crahses: Number of found crashes
  - Origin: An original bug-triggering seed
  - Origin length: Length of `Origin`
  - Minimized: Minimized `Origin`
  - Minimized length: Length of `Minimized`
- Coverage
  - File: Name of target program file
  - Function hit ratio: Function hit ratio from coverage module (chapter 1)
  - Branch hit ratio: Branch hit ratio from coverage module (chapter 1)
  - Line hit ratio: Line hit ratio from coverage module (chapter 1)
- Seed
  - Seeds: Number of generated seeds
  - Current seed: A current seed to be fed
  - Current seed score: Score of `Current seed`
  - Visited edges: Visited number of basic blocks
  - Next seed: A seed to be fed at next round
  - New paths: Number of found new edges.
- Stdout
  - Display the stdout of target program

Note that the Coverage panel will not display live numbers unless the coverage instrumentation and its corresponding runtime are actively involved. The following function defines five distinct data structures, one for each panel.

```
1  // fuzzer/src/campaign.rs Rust
2
3  #[derive(Debug)]
4  pub struct CrashInfo {
5      pub crashes: usize,
6      pub origin: Seed,
7      pub minimized: Seed,
8  }
9  impl CrashInfo {
10     pub fn new(crashes: usize, origin: Seed, minimized: Seed) ->
        Self {
11         return Self {
12             crashes,
13             origin,
14             minimized,
15         };
16     }
17 }
18
19 #[derive(Debug)]
20 pub struct FuzzerMetadata {
21     pub fuzz_cnt: u64,
22     pub fuzz_input_typ: FuzzInput,
23     pub timeout: Duration,
24     pub target_elapsed_time: Duration,
25     pub total_elapsed_time: Duration,
26 }
27 impl FuzzerMetadata {
28     pub fn new(
29         fuzz_cnt: u64,
30         fuzz_input_typ: FuzzInput,
31         timeout: Duration,
32         target_elapsed_time: Duration,
33         total_elapsed_time: Duration,
34     ) -> Self {
35         Self {
36             fuzz_cnt,
37             fuzz_input_typ,
```

```
38         timeout,
39         target_elpased_time,
40         total_elpased_time,
41     }
42 }
43 }
44
45 #[derive(Debug)]
46 pub struct FuzzerSeed {
47     pub seeds: usize,
48     pub cur_seed: Seed,
49     pub next_seed: Seed,
50     pub visit_edges: u64,
51     pub new_paths: usize,
52 }
53 impl FuzzerSeed {
54     pub fn new(
55         seeds: usize,
56         cur_seed: Seed,
57         next_seed: Seed,
58         visit_edges: u64,
59         new_paths: usize,
60     ) -> Self {
61         Self {
62             seeds,
63             cur_seed,
64             next_seed,
65             visit_edges,
66             new_paths,
67         }
68     }
69 }
70
71 #[derive(Debug)]
72 pub enum FuzzShot {
73     ProgramOutput(String),
74     Coverage(Vec<CovReport>),
75     Crash(CrashInfo),
76     Metadata(FuzzerMetadata),
77     SeedInfo(FuzzerSeed),
78     Terminated,
```

```
79  }
80
81  #[derive(Debug)]
82  pub struct FuzzingCampaign {
83      pub program_output: Vec<String>,
84      pub coverage: Vec<CovReport>,
85      pub crash: CrashInfo,
86      pub metadata: FuzzerMetadata,
87      pub seed_info: FuzzerSeed,
88  }
89  impl FuzzingCampaign {
90      pub fn new(fuzz_input: FuzzInput) -> Self {
91          Self {
92              program_output: vec!["".to_string()],
93              coverage: vec![CovReport {
94                  file: "".to_string(),
95                  funcs_hit_ratio: "".to_string(),
96                  uncovered_funcs: "".to_string(),
97                  brs_hit_ratio: "".to_string(),
98                  uncovered_brs: "".to_string(),
99                  lines_hit_ratio: "".to_string(),
100                 uncovered_lines: "".to_string(),
101             }],
102             crash: CrashInfo::new(0, Seed::new(vec![], 0),
103                                   Seed::new(vec![], 0)),
104             metadata: FuzzerMetadata::new(
105                 0,
106                 fuzz_input,
107                 Duration::default(),
108                 Duration::default(),
109                 Duration::default(),
110             ),
111             seed_info: FuzzerSeed::new(0, Seed::new(vec![], 0),
112                                         Seed::new(vec![], 0), 0, 0),
113         }
114     }
115     pub fn add_program_output(&mut self, pgm_output: String) {
116         self.program_output.push(pgm_output);
117     }
118     pub fn set_coverage(&mut self, coverage: Vec<CovReport>) {
119         if !coverage.is_empty() {
120             self.coverage = coverage;
121         }
122     }
123 }
```



```

119     }
120 }
121 pub fn set_crash(&mut self, crash: CrashInfo) {
122     self.crash = crash;
123 }
124 pub fn set_metadata(&mut self, metadata: FuzzerMetadata) {
125     self.metadata = metadata;
126 }
127 pub fn set_seed_info(&mut self, seed_info: FuzzerSeed) {
128     self.seed_info = seed_info;
129 }
130 }

```

The fuzzer sends these objects to the UI application, which then updates the live interface. Instead of presenting the entire code here, as it exceeds 300 lines, we'll focus only on the core logic.

```

1 // fuzzer/src/ui.rs Rust
2 pub fn run_ui(rx: mpsc::Receiver<FuzzShot>) -> io::Result<()> {
3     enable_raw_mode()?;
4     let mut stdout = stdout();
5     execute!(stdout, EnterAlternateScreen)?;
6     let backend = CrosstermBackend::new(stdout);
7     let mut terminal = Terminal::new(backend)?;
8
9     let res = run_app(&mut terminal, rx)?;
10
11     disable_raw_mode()?;
12     execute!(terminal.backend_mut(), LeaveAlternateScreen)?;
13     terminal.show_cursor()?;
14     println!("{}", res);
15     Ok(())
16 }
17
18 fn run_app<B: ratatui::backend::Backend>(
19     terminal: &mut Terminal<B>,
20     rx: mpsc::Receiver<FuzzShot>,
21 ) -> io::Result<String> {
22     let mut fuzz_result = FuzzingCampaign::new(FuzzInput::Stdin);
23     loop {
24         match rx.recv() {
25             Ok(result) => match result {

```

```

26      FuzzShot::ProgramOutput(o) => {
27          fuzz_result.add_program_output(o);
28      }
29      FuzzShot::Coverage(cov) => {
30          fuzz_result.set_coverage(cov);
31      }
32      FuzzShot::Crash(crash) => {
33          fuzz_result.set_crash(crash);
34      }
35      FuzzShot::Metadata(md) => {
36          fuzz_result.set_metadata(md);
37      }
38      FuzzShot::SeedInfo(seed) => {
39          fuzz_result.set_seed_info(seed);
40      }
41      FuzzShot::Terminated => {
42          return Ok("Fuzzer terminated".to_string());
43      }
44  },
45  Err(err) => {
46      panic!("{:?}", err);
47  }
48  }
49
50  terminal.draw(|frame| render_ui(frame, &mut fuzz_result))?;
51
52  let timeout = Duration::from_millis(10);
53  if event::poll(timeout)? {
54      if let Event::Key(key) = event::read()? {
55          if KeyCode::Char('q') == key.code {
56              return Ok("'Q' entered".to_string());
57          }
58      }
59  }
60  }
61  }

```

The loop at line 23 receives live fuzzing data from the fuzzer, which then updates the UI at line 50. If the user presses 'q', the fuzzer terminates, the quit handler is presented at lines 52-60.

### 5.4.3.8 Fuzz run

With all implementations complete, we can now run the fuzzer in two steps: instrumentation and fuzzer execution.

Here's the example target program we'll be fuzzing:

```
1 // target.c
2 #include <stdio.h>
3
4 #define FUZZ_LEN 10
5 #define ABNORMAL_EXIT 100
6
7 int main(int argc, char* argv[]) {
8     char input[FUZZ_LEN];
9     fgets(input, sizeof(input), stdin);
10    printf("target run\n");
11
12    if (input[0] == 'a') {
13        if (input[1] == 'c' || input[1] == 'd' || input[1] == 'e' ||
14            input[1] == 'f') { // buggy path
15            // input[FUZZ_LEN] = 'a'; // you can use to trigger a
16                // explicit bug using OOB
17            return ABNORMAL_EXIT;
18        }
19    }
20    return 0;
21 }
```

This program accepts user input via the `stdin` file descriptor. We've intentionally introduced a bug at either line 14 or 15 (you can choose which one). This bug triggers if the input seed's first byte is 'a', and the `stdin` input also contains 'c', 'd', 'e', or 'f' at index 1.

Since we've implemented ASan, we can trigger an explicit OOB access to cause a crash. Alternatively, we can make the program return any value other than 0 or 1 to signal an error.

These are the compile commands:

```
1 clang -g -O0 -S -emit-llvm target.c -o target.ll
2 ./instrument -i target.ll -o itarget.ll -m all
3 clang itarget.ll -L. -lcoverage_runtime -lfuzzer_runtime -
4     lasan_runtime
5 // the binary 'a.out' will be generated.
```

You'll need to generate your initial seed files within the `init-seeds` directory. For demonstration, we've simply added three seed files: `a.txt`, `b.txt`, and `c.txt`.

```
1 init-seeds/a.txt: aaba
2 init-seeds/b.txt: aaab
```

```
3 init-seeds/c.txt: abaa
```

Now we can run the fuzzer with the following command: `./fuzzer -p ./a.out -c init-seeds -t stdin`.

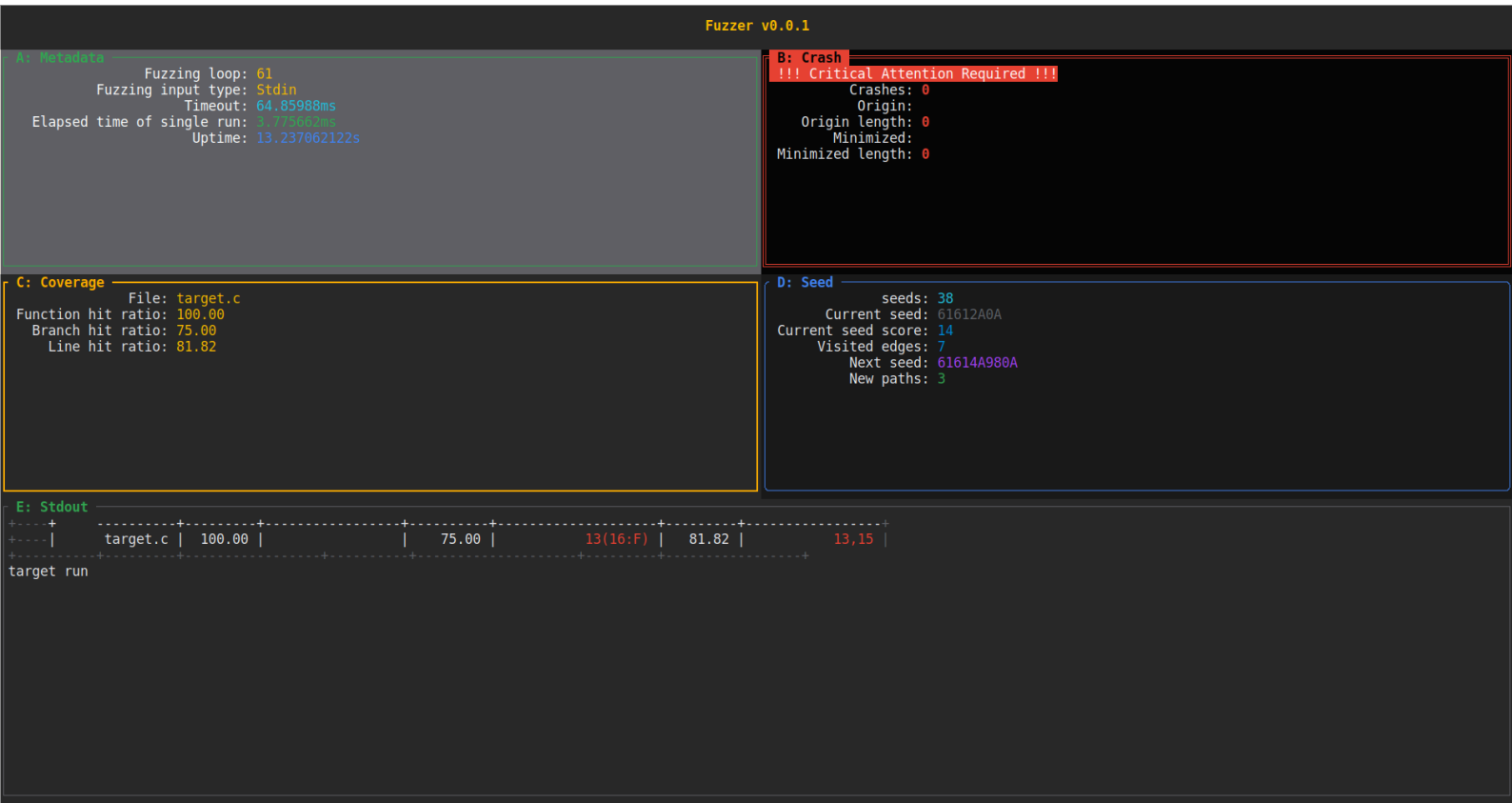


Figure 5.6: Fuzzer Execution

Crashes will be stored in the `crashes` directory, which is generated automatically.

```
1 crashes/1.crash: ae
2 crashes/2.crash: ad
3 crashes/3.crash: ac
```

The fuzzer successfully finds a buggy path for this example target program in under a minute on average. However, if you change the condition at line 13 to `input[1] == 'c'`, it'll take more time to trigger the bug. This is because the probability of generating a seed that satisfies this new condition is significantly lower (.because previously, the value at index 1 could be 'c', 'd', 'e', or 'f'. Now, it specifically requires 'c', ignoring 'd', 'e', and 'f').

5.4.4 Summary

We fully implemented AFL-like fuzzing campaign by designing and implementing instrumentation module, runtime, and fuzzer.

#### 5.4.4.1 Homework

**Exercise 5.1** We suggest the following topics for further study and exploration:

- **Topic #1: Comparing Vanilla Fuzzer vs. Fork-Server Based Fuzzer Performance**
  - You can experiment with a vanilla fuzzer that doesn't utilize the fork-server mechanism (i.e., `execve()` is called at every run), loading the target program at each fuzzing round. This will allow you to learn how much a fork-server based fuzzer contributes to boosting fuzzing speed.
- **Topic #2: Smart Mutation Strategy**
  - Currently, all mutation methods are selected randomly. However, you could assign weights to each method, giving higher probabilities to certain techniques under specific circumstances. By employing a heuristic approach to develop a smarter strategy manager, you could then compare its performance to a random selection and evaluate if it contributes to finding bugs more quickly.

Numerous highly-tuned fuzzers are regularly submitted to and published in Computer Science conferences. You are encouraged to explore these publications to discover interesting topics and study current trends and challenges in the field. We hope this area inspires you to pursue rewarding research and coding endeavors.



## Part V

# Symbolic Execution

### Introduction to Symbolic

<b>6</b>	<b>Execution . . . . .</b>	<b>129</b>
6.1	Problem . . . . .	129
6.2	Working Example . . . . .	130
6.3	Design . . . . .	131
6.4	Implementation . . . . .	134





# “How to satisfy the branch condition to proceed?”

```
if (a + 5 == 42) {  
    ... // Block 1  
} else {  
    ... // Block 2  
}
```



```
Block 1: a == 37  
Block 2: a != 37
```

## 6. Introduction to Symbolic Execution

### 6.1 Problem

In the previous chapter, we implemented a coverage-based fuzzer. The goal of this fuzzer is to determine the fuzzing direction through its coverage footprint. It's important to understand that while it uses coverage to guide its decisions, it doesn't inherently increase coverage itself. To truly boost the coverage ratio, a fuzzer needs to be aware of how to satisfy all branches so it can enter them, thereby increasing code coverage. For example, the branch condition below would be very difficult to satisfy via mutation alone.

```
if (input() ==  
1 "0x9e6d3c81b42a7f124e5f9a51c23df61d5a97a4c3f7e9b881b3e4d2a6b92cbbaf")  
{  
2 ...  
3 } else {  
4 ...  
5 }
```

Most coverage-based fuzzers will never reach line 2, consistently taking only the else branch in every run. This is because they lack any hint or mechanism for generating inputs that could satisfy the conditions needed to reach other branches.

Let's consider a simpler example.

```
1 if (input() > 100) {  
2 ...  
3 } else {  
4 ...  
5 }
```

To take the true branch, the `input()` value must be greater than 100. A coverage-based fuzzer might discover a value exceeding 100 through mutation, but this is purely a matter

of random chance. This lack of awareness regarding branch conditions is a significant hurdle when trying to generate appropriate seed values for specific target branches.

Symbolic execution addresses this challenge by providing the ability to find a solution that satisfies a target branch. In this book, we won't be implementing a hybrid fuzzer that combines coverage-based techniques (from Chapter 3) with symbolic execution (from Chapter 4). This will be detailed in the homework exercises at the end of this chapter.

## 6.2 Working Example

To interact with symbolic execution, the programmer must symbolize specific variables. The values for these variables are then provided by the symbolic execution runtime. The following code is an example of symbolization.

```
1 // symbolic-example.c
2 #include <stdio.h>
3
4 extern void __make_symbolic();
5
6 int main() {
7     int i,j;
8     __make_symbolic(sizeof(int), &i);
9     __make_symbolic(sizeof(int), &j);
10    printf("%d, %d\n", i,j);
11    if (i == 0) {
12        if (j == 123214125) {
13            printf("S1\n");
14        } else {
15            printf("S2\n");
16        }
17    } else {
18        if (j != 88148128) {
19            printf("S3\n");
20        } else {
21            printf("S4\n");
22            return 123;
23        }
24    }
25 }
```

The example target program contains six branches, and it evaluates and compares the symbolic variables `i` and `j`. The symbolic execution runtime provides the values for these variables. Here are the compilation commands. Note that we're using the `-m all` option, which enables all instrumentation modules.

```
1 clang -g -O0 -S -emit-llvm symbolic-example.c -o symbolic-example.ll
2 ./instrument -i symbolic-example.ll -o isymbolic-example.ll -m all
```

```
3 clang isymbolic-example.ll -L. -lcovrage_runtime -lasan_runtime -
   lfuzzer_runtime -lsymbolic_runtime
```

Let's run the generated binary.

```
1  > LD_LIBRARY_PATH=./ ./a.out
2  1, 123214126
3  S3
4  ... (coverage report)
5
6  > LD_LIBRARY_PATH=./ ./a.out
7  1, 88148128
8  S4
9  ... (coverage report)
10
11 > LD_LIBRARY_PATH=./ ./a.out
12 0, 123214125
13 S1
14 ... (coverage report)
15
16 > LD_LIBRARY_PATH=./ ./a.out
17 0, 88148128
18 S2
19 ... (coverage report)
```

Even after multiple runs, all branches are covered. However, with a coverage-based fuzzer, generating specific values like `i == 0 && j = 123214125` (printing “S1”) or `i == 1 && j = 88148128` (printing “S4”) would be extremely difficult. While the final binary can be integrated with our coverage-based fuzzer (`./fuzzer -p ./a.out -c init-seeds -t stdin`), the fuzzer will essentially just repeat the execution process. This means the symbolic values will still be provided by the symbolic execution runtime, and the fuzzer itself won't be aware of the symbolic execution process. A detailed look at the ideal integration will be covered in the homework section.

## 6.3 Design

Symbolic execution runtime manages condition branches where symbolic variables are used. Specifically it manages “state” which values satisfies or not each branch conditions.

### 6.3.1 State

```
1 // symbolic-example.c
2 #include <stdio.h>
3
4 extern void __make_symbolic();
5
```

C

```
6  int main() {
7      int i,j;
8      __make_symbolic(sizeof(int), &i);
9      __make_symbolic(sizeof(int), &j);
10     printf("%d, %d\n", i,j);
11     // state 0
12     if (i == 0) { // state 1
13         if (j == 123214125) { // state 2
14             printf("S1\n");
15         } else { // state 3
16             printf("S2\n");
17         }
18     } else { // state4
19         if (j != 88148128) { // state 5
20             printf("S3\n");
21         } else { // state 6
22             printf("S4\n");
23             return 123;
24         }
25     }
26 }
```

In our example, there are a total of six states. A new state is created whenever a branch condition incorporates a symbolic value. Let's trace this process starting from line 7.

Lines 7 declares two integer variables, *i* and *j*. Lines 8 and 9 indicate that *i* and *j* are symbolic variables, meaning their concrete values will be supplied by the symbolic execution runtime.

At line 12, the symbolic variable *i* is compared. At this point, a new state, State 1, is created. Additionally, another state, State 2 (line 18), is also generated. When a new state is required due to a branch, two states are created through a process called **forking**.

A fork involves copying the current state and only changing its predicate by negating it. For example, when creating State 2 and State 3, they will have identical content except for their predicates. State 2 will have the “equal” predicate, while State 3 will have the “not equal” predicate.

At line 13, the symbolic variable *j* is used in a branch condition, leading to the creation of State 3 and State 4 (line 15). Finally, line 19 creates State 5 and State 6 (line 22). The figure below illustrates each of these state nodes.

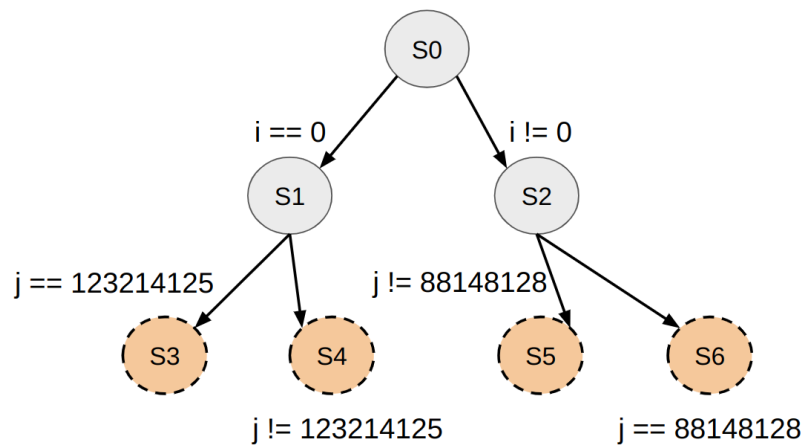


Figure 6.1: Symbolic State

This gives us a total of six states, excluding the initial state, ‘S0’.

### 6.3.2 Static Analysis vs. Dynamic Management

I believe there are two main approaches to managing states: static analysis or dynamic analysis.

If you delegate state creation to runtime, dynamic analysis comes into play. When the target program starts, nothing special happens initially. However, when a branch involving a symbolic variable is encountered, the system detects the relevant instruction (like `icmp`) and dynamically creates a new state. This approach is typically used when the target program is assumed to be, or limited to, a binary.

When source code is available, we can opt for static analysis. This allows us to identify how many branch conditions involve symbolic variables at compile time (during instrumentation). This way, potential states can be calculated upfront. However, static analysis can’t pre-calculate all possible states. This challenge can be mitigated with concolic testing [25] (combining concrete and symbolic execution); by executing concrete inputs and using the feedback, it can prune useless states. We’ll detail this in the homework section.

In this book, we’re opting for a static analysis approach to create all explicit states (as demonstrated in the limited example above). For instance, our implementation achieves 100% coverage on the working example. However, if a conditional branch requires an unknown state—a combination of existing states that can’t be resolved during static analysis—it introduces the problems of state pruning and state explosion. We’ll delve into these challenges further in the homework section.

Returning to our specific scope, we instrument the encoding for all explicit states. Our runtime then resolves the resulting equations using Z3 [26], a powerful theorem prover.

### 6.3.3 Interface (Harness)

In our example, the `__make_symbolic()` functions at lines 8 and 9 serve as the interface between the target program and the symbolic runtime. During each run, the symbolic runtime uses Z3 to supply resolved values for these variables. Therefore, when perform-

ing symbolic testing on any program, you'll need to select your symbolic variables and initialize them using this function, essentially creating a testing harness.

## 6.4 Implementation

As always, our implementation consists of two parts: the instrumentation module and the runtime.

### 6.4.1 Instrumentation

Our instrumentation strategy is divided into four steps:

- **Collecting Symbolic Variables:** We capture symbolic variables by identifying calls to the `__make_symbolic()` initialization function.
- **State Creation:** By iterating through basic blocks, we look for `icmp` instructions. If a symbolic variable is involved, we create a new state by forking
- **State Initialization via Runtime:** Each state node is encoded and passed to a runtime initialization function. Only the leaf state nodes are serialized and passed into this function.
- **Instrumentation:** The serialized states are instrumented within a constructor function, and all constraints from the target program are initialized in the runtime

Let's begin by looking at the data structure.

```
1 // instrument/src/symbolic.rs Rust
2
3 pub const VAR_KIND: i8 = 0;
4 pub const CONST_KIND: i8 = 1;
5
6 pub const PREDICATE_EQ: i8 = 0;
7 pub const PREDICATE_NE: i8 = 1;
8 pub const PREDICATE_SLT: i8 = 2;
9 pub const PREDICATE_SLE: i8 = 3;
10 pub const PREDICATE_SGT: i8 = 4;
11 pub const PREDICATE_SGE: i8 = 5;
12
13 #[derive(Debug, Clone)]
14 enum Operand<'ctx> {
15     Var(PointerValue<'ctx>),
16     Const(i64),
17 }
18 impl<'ctx> Operand<'ctx> {
19     fn is_const(&self) -> bool {
20         if let Self::Const(_) = self {
21             return true;
22         } else {
23             return false;
```

```

24     }
25 }
26 }
27
28 #[derive(Debug, Clone)]
29 struct Constraint<'ctx> {
30     left_operand: Operand<'ctx>,
31     right_operand: Operand<'ctx>,
32     predicate: IntPredicate,
33 }
34
35 impl<'ctx> Constraint<'ctx> {
36     fn new(
37         left_operand: Operand<'ctx>,
38         right_operand: Operand<'ctx>,
39         predicate: IntPredicate,
40     ) -> Self {
41         Self {
42             left_operand,
43             right_operand,
44             predicate,
45         }
46     }
47
48     fn predicate_to_i8(&self) -> i8 {
49         match &self.predicate {
50             IntPredicate::EQ => PREDICATE_EQ,
51             IntPredicate::NE => PREDICATE_NE,
52             IntPredicate::SLT => PREDICATE_SLT,
53             IntPredicate::SLE => PREDICATE_SLE,
54             IntPredicate::SGT => PREDICATE_SGT,
55             IntPredicate::SGE => PREDICATE_SGE,
56             _ => unreachable!("unsigned operation is not
57                 supported"),
58         }
59
60     fn negate(&self) -> Self {
61         let predicate = match self.predicate {
62             IntPredicate::EQ => IntPredicate::NE,
63             IntPredicate::NE => IntPredicate::EQ,

```

```

64         IntPredicate::SLT => IntPredicate::SGE,
65         IntPredicate::SLE => IntPredicate::SGT,
66         IntPredicate::SGT => IntPredicate::SLE,
67         IntPredicate::SGE => IntPredicate::SLT,
68         _ => unreachable!("unsigned operation is not
        supported"),
69     };
70     Self::new(
71         self.left_operand.clone(),
72         self.right_operand.clone(),
73         predicate,
74     )
75 }
76 }

```

The `Operand` struct represent whether the operand is symbolic variable or constant value. In the working example, the symbolic variables `i` and `j` correspond to `Var` and constant values correspond to `Const` type such as 123214125 and 88148128.

The `Constraint` struct represents a branch condition. Branch condition contains left and right operand and predictor. For example, `icmp eq i32 %1, 555` includes left operand `%1(Var)`, right operand(`555(Const)`) and predictor (`IntPredicate::EQ`). We consider only six type predictor, `==`, `!=`, `<`, `<=`, `>`, and `>=`. You can include other predictors for extension.

The function `predicate_to_i8()` is used when a state is serialized. The `negate()` function

The `negate()` function negate the predictor and other properties are remain.

The `Operand` struct tells us if an operand is a symbolic variable or a constant value. In our example, symbolic variables like `i` and `j` would be represented by the `Var` type, while constant values such as 123214125 and 88148128 would be `Const` types.

The `Constraint` struct represents a branch condition. It holds a left operand, a right operand, and a predicate. For instance, an instruction like `icmp eq i32 %1, 555` would have `%1(Var)` as its left operand, `555(Const)` as its right operand, and `IntPredicate::EQ` as its predicate. We're focusing on six types of predicates: `==`, `!=`, `<`, `<=`, `>`, and `>=`. You can extend this to include others if needed.

The `predicate_to_i8()` function is used when a state is serialized.


The `negate()` function, on the other hand, simply negates the predicate while keeping all other properties unchanged.

Now we introduce state structure:

```

1  // instrument/src/symbolic.rs
2
3  #[derive(Debug, Clone)]
4  struct State<'ctx> {
5      id: i64,
6      path_constraints: Vec<Constraint<'ctx>>,
7      is_leaf: bool,

```

 Rust



```

8   }
9
10  impl<'ctx> State<'ctx> {
11      fn new() -> Self {
12          Self {
13              id: 0,
14              path_constraints: vec![],
15              is_leaf: false,
16          }
17      }
18
19      fn fork(&self, constraint: Constraint<'ctx>) -> (Self, Self) {
20          let mut id = CONSTRAINT_ID.lock().unwrap();
21          let (tid, fid) = (*id + 1, *id + 2);
22          *id += 2;
23          let (mut copied_state_t, mut copied_state_f) =
24              (self.clone(), self.clone());
25          copied_state_t.path_constraints.push(constraint.clone());
26          copied_state_f.path_constraints.push(constraint.negate());
27          copied_state_t.id = tid;
28          copied_state_f.id = fid;
29          (copied_state_t, copied_state_f)
30      }
31  }

```

Each state is assigned a unique ID, which serves as its identifier. The `path_constraints` field holds a set of branch conditions. It's crucial that a branch condition within a nested structure preserves the constraints from its previous state. Let's look at an example to understand why preservation is required:

```

1  if (a == 1) {
2      if (b == 1) {
3          ...
4      }
5  }

```

To reach line 3, we must satisfy two conditions: `a == 1` and `b == 1`. If either of these isn't met, line 3 will never execute. This highlights why we need to be aware of the current state before forking a new one. In the example above, at line 2, we have one constraint: `a == 1`. When we create a new state, we append the new branch condition to the current state, resulting in `a == 1 && b == 1`. This is precisely why `path_constraints` is implemented as a list (or vector) type.

To create a fork, we generate two new IDs and clone the current state twice. One clone represents the true branch, while the other represents the false branch. For the

false branch, we simply negate the predicate. The implementation is straightforward; you can examine the `fork()` function for details.

Now, let's dive into the main logic of our instrumentation. We'll focus on the core components here; you can find the full implementation in the repository.

```

1  // instrument/src/symbolic.rs Rust
2
3  #[derive(Default)]
4  pub struct SymbolicModule {}
5
6  fn build_sym_ptrs<'ctx>(
7      context: &'ctx Context,
8      module: &Module<'ctx>,
9      builder: &Builder<'ctx>,
10 ) -> Result<Vec<PointerValue<'ctx>>> {
11     let mut symbolic_ptr_inits = HashSet::new();
12     let funcs: Vec<_> = module.get_functions().collect();
13     for func in funcs {
14         if can_skip_instrument(&func) {
15             continue;
16         }
17         for basic_blk in func.get_basic_blocks() {
18             for instr in basic_blk.get_instructions() {
19                 if instr.get_opcode() == InstructionOpcode::Call &&
20                    instr.get_num_operands() == 3 {
21                     let func_ptr =
22                         instr.get_operand(2).unwrap().left().unwrap();
23                     let func_str = cstr_to_str(func_ptr.get_name());
24                     if func_str == SYMBOLIC_MAKE_VAR { //
25                         __make_symbolic()
26                     }
27                     let var_ptr = instr
28                         .get_operand(1)
29                         .unwrap()
30                         .left()
31                         .unwrap()
32                         .into_pointer_value();
33                     symbolic_ptr_inits.insert(var_ptr);
34
35                     // install <address, pointer value> mapping
36                     builder.position_before(&instr);
37                     build_sym_make_prep(context, module,
38                                         builder, var_ptr)?;
39                 }
40             }
41         }
42     }
43     Ok(symbolic_ptr_inits)
44 }
```

```

35         }
36     }
37 }
38 }
39 Ok(symbolic_ptr_inits.iter().cloned().collect())
40 }

```

The function `build_sym_ptr()` collects all symbolic variables by looking the function `__make_symbolic()`.

```

1  // instrument/src/symbolic.rs Rust
2
3  fn get_sym_ptr_operand<'ctx>(<
4      sym_ptrs: &Vec<PointerValue<'ctx>>,
5      op: BasicValueEnum<'ctx>,
6  ) -> Option<Operand<'ctx>> {
7      if op.is_int_value() {
8          let op = op.into_int_value();
9          if op.is_const() {
10             return
11             Some(Operand::Const(op.get_sign_extended_constant().unwrap()));
12         } else {
13             let mut instr = op.as_instruction().unwrap();
14             loop {
15                 match instr.get_opcode() {
16                     InstructionOpcode::Load => {
17                         let ptr = instr.get_operand(0).
18                             unwrap().
19                             left().unwrap();
20                         for sym_ptr in sym_ptrs {
21                             if *sym_ptr == ptr {
22                                 return Some(Operand::Var(*sym_ptr));
23                             }
24                         }
25                         return None;
26                     }
27                     _ => {
28                         if let Some(prev_instr) =
29                             instr.get_previous_instruction() {
30                             instr = prev_instr;
31                         } else {
32                             return None;

```

```

32             }
33         }
34     }
35 }
36 }
37 }
38     None
39 }

```

The `get_sym_ptr_operand()` function is designed to retrieve symbolic variables within branch conditions.

- If a branch condition's operand is a constant integer, the function returns the constant, wrapped appropriately.
- If the operand is a symbolic variable, it returns the symbolic variable, also wrapped.
- If the operand is neither (i.e., not a symbolic variable), `None` is returned, indicating that no symbolic variable is involved in that particular branch condition.

The code below illustrates the state initialization instrumentation. We split into two parts because of the page limit.

```

1  // [1-1] Rust
2  // instrument/src/symbolic.rs
3
4  impl InstrumentModule for SymbolicModule {
5      fn instrument<'ctx>(<
6          &self,
7          context: &'ctx Context,
8          module: &Module<'ctx>,
9          builder: &Builder<'ctx>,
10     ) -> Result<()> {
11         let sym_ptrs = build_sym_ptrs(context, module, builder)?;
12         let mut serialized = vec![];
13         let funcs: Vec<_> = module.get_functions().collect();
14         for func in funcs {
15             let mut states: HashMap<i64, State> = HashMap::new();
16             let mut bb_id = HashMap::new();
17             if let Some(first_bb) = func.get_first_basic_block() {
18                 let first_bb_addr = first_bb.as_mut_ptr() as usize;
19                 bb_id.insert(first_bb_addr, 0);
20                 states.insert(0, State::new());
21             }
22             for basic_blk in func.get_basic_blocks() {
23                 for instr in basic_blk.get_instructions() {
24                     .... // refer the code snippet below [1-2]

```

```

25         }
26         instrumented_blks.insert(basic_blk);
27     }
28     if states.len() > 1 {
29         serialized.push(serialize_constraints(&states));
30     }
31 }
32 if !serialized.is_empty() {
33     let constructor = build_symbolic_init(context, module,
34     builder, serialized)?;
35     build_ctros(context, module, constructor)?;
36 }
37 // Verify instrumented IRs
38 module_verify(module)
39 }
40 }

```

Our goal is to instrument serialized constraints so the runtime can recognize which constraints exist.

We begin by identifying the symbolic variables. We insert State 0 at the first basic block of each function, meaning our current implementation is limited to intra-procedural analysis. State 0 doesn't contain any constraints itself.

The core logic resides in the loop from lines 22 to 27, with detailed code and descriptions provided below.

Once states are created, they're serialized and added to the serialized vector. After iterating through all states, these serialized constraints are then instrumented within a constructor function.

```

1  // instrument/src/symbolic.rs Rust
2  // [1-2]
3
4  match instr.get_opcode() {
5      InstructionOpcode::Br => {
6          // if the current basic block has a state and the terminator is
6          // not conditional, we mark them as leaf
7          if !instr.is_conditional() {
8              let bb_addr = basic_blk.as_mut_ptr() as usize;
9              if let Some(id) = bb_id.get(&bb_addr) {
10                 if let Some(state) = states.get_mut(&id) {
11                     state.is_leaf = true;
12                 }
13             }
14         }

```



```

55         states.insert(state_f.id, state_f);
56     }
57 } else {
58     let state = states.get(&0).unwrap();
59     let (state_t, state_f) =
        state.fork(constraint.clone());
60     bb_id.insert(tbr.as_mut_ptr() as usize, state_t.id);
61     bb_id.insert(fbr.as_mut_ptr() as usize, state_f.id);
62     states.insert(state_t.id, state_t);
63     states.insert(state_f.id, state_f);
64 }
65 }
66 _ => {}
67 }
68 }
69 }
70 }
71 _ => {}
72 }

```

The code above represents the main handler for collecting states. It identifies conditional branch instructions and then determines if their operands correspond to symbolic variables. If a symbolic variable is found, it creates a `Constraint` struct and generates two new states by forking. As we mentioned earlier, to correctly handle nested branches, the previous state is queried using its basic block address and state ID.

When a new state is added, the addresses for both the true and false branches are stored, along with the IDs of the two newly generated states. This ensures that the current state information propagates to future states, allowing them to correctly determine which states need to be considered.

The following function serializes the set of constraints.

```

1  // instrument/src/symbolic.rs Rust
2
3  #[derive(Debug)]
4  pub struct ConstraintSerialized {
5      pub id: i64,
6      pub left_operand_kind: i8,
7      pub left_operand_val: i64,
8      pub right_operand_kind: i8,
9      pub right_operand_val: i64,
10     pub predicate: i8,
11 }
12
13 impl ConstraintSerialized {

```

```

14     pub fn new(
15         id: i64,
16         left_operand_kind: i8,
17         left_operand_val: i64,
18         right_operand_kind: i8,
19         right_operand_val: i64,
20         predicate: i8,
21     ) -> Self {
22         Self {
23             id,
24             left_operand_kind,
25             left_operand_val,
26             right_operand_kind,
27             right_operand_val,
28             predicate,
29         }
30     }
31 }
32
33 fn serailize_constraints<'ctx>(states: &HashMap<i64, State>) ->
Vec<ConstraintSerialized> {
34     let mut constraints = vec![];
35     for (_, state) in states {
36         if state.is_leaf {
37             for constraint in &state.path_constraints {
38                 let (left_kind, left_op) = match
constraint.left_operand {
39                     Operand::Var(var) => (VAR_KIND,
var.as_value_ref() as i64),
40                     Operand::Const(v) => (CONST_KIND, v),
41                 };
42                 let (right_kind, right_op) = match
constraint.right_operand {
43                     Operand::Var(var) => (VAR_KIND,
var.as_value_ref() as i64),
44                     Operand::Const(v) => (CONST_KIND, v),
45                 };
46                 let predicate = constraint.predicate_to_i8();
47                 constraints.push(ConstraintSerialized::new(
48                     state.id, left_kind, left_op, right_kind,
right_op, predicate,
49                 ));
50             }

```



```

51     }
52 }
53 constraints
54 }

```

The only leaf nodes are instrumented as follows in the given working example.

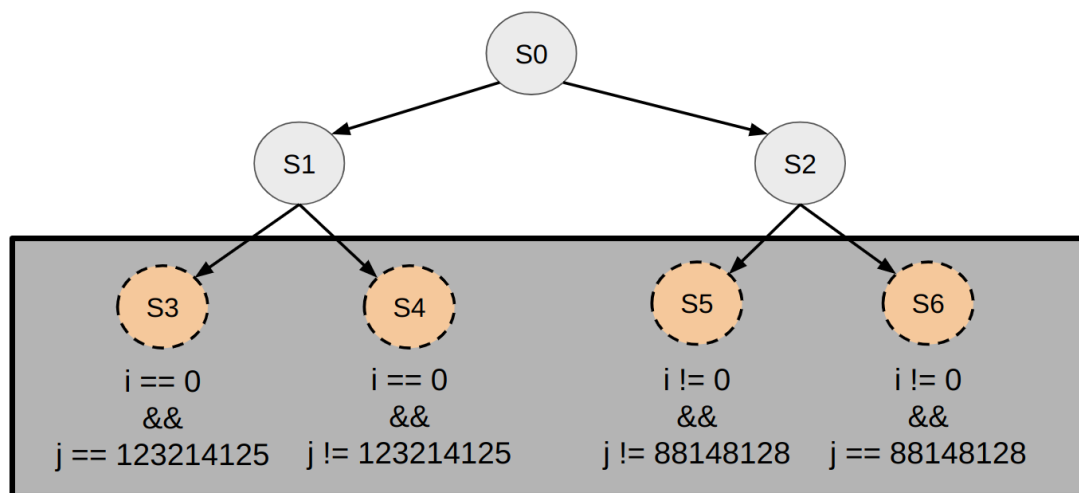


Figure 6.2: Leaf State Nodes

Finally, each serialized states are instrumented as follows:

```

1  define i32 @main() #0 !dbg !27 { LLVM-IR
2    ...
3    call void @__symbolic_make_prepare(ptr %3, i64 106096555351408), !
  dbg !36
4    call void (i64, ptr, ...) @__make_symbolic(i64 noundef 4, ptr
  noundef %3), !dbg !36
5    call void @__symbolic_make_prepare(ptr %4, i64 106096555351136), !
  dbg !37
6    call void (i64, ptr, ...) @__make_symbolic(i64 noundef 4, ptr
  noundef %4), !dbg !37
7    ...
8    ...
9
10 declare void @__symbolic_module_add_sym(i64, i8, i64, i8, i64, i8)
11
12 define void @__symbolic_init(i64 %0, i8 %1, i64 %2, i8 %3, i64 %4,
  i8 %5) {
13   __symbolic_init_entry:

```

```

14  call void @_symbolic_module_add_sym(i64 5, i8 0, i64
    106096555351408, i8 1, i64 0, i8 1)
15  call void @_symbolic_module_add_sym(i64 5, i8 0, i64
    106096555351136, i8 1, i64 88148128, i8 1)
16  call void @_symbolic_module_add_sym(i64 4, i8 0, i64
    106096555351408, i8 1, i64 0, i8 0)
17  call void @_symbolic_module_add_sym(i64 4, i8 0, i64
    106096555351136, i8 1, i64 123214125, i8 1)
18  call void @_symbolic_module_add_sym(i64 6, i8 0, i64
    106096555351408, i8 1, i64 0, i8 1)
19  call void @_symbolic_module_add_sym(i64 6, i8 0, i64
    106096555351136, i8 1, i64 88148128, i8 0)
20  call void @_symbolic_module_add_sym(i64 3, i8 0, i64
    106096555351408, i8 1, i64 0, i8 0)
21  call void @_symbolic_module_add_sym(i64 3, i8 0, i64
    106096555351136, i8 1, i64 123214125, i8 0)
22  ret void
23  }

```

Since only the leaf nodes are instrumented, there are a total of eight states, resulting from the forking of four distinct states.

### 6.4.2 Runtime

The core of our runtime implementation involves interacting with the Z3 solver.

Our first task is to add constraints for each state ID.

```

1  symbolic_runtime/src/runtime.rs Rust
2
3  #[no_mangle]
4  pub extern "C" fn __symbolic_module_add_sym(
5      id: i64,
6      left_operand_kind: i8,
7      left_operand_val: i64,
8      right_operand_kind: i8,
9      right_operand_val: i64,
10     predicate: i8,
11 ) {
12     CONSTRAINTS.write().unwrap().add_constraint(
13         id,
14         left_operand_kind,
15         left_operand_val,
16         right_operand_kind,
17         right_operand_val,
18         predicate,
19     );

```

```
20 }
```

The global singleton CONSTRAINTS holds all constraints associated with each state ID. This means that all constraints within a given ID are considered collectively to find a solution.

```
1  // symbolic_runtime/src/symbolic.rs  Rust
2
3  lazy_static::lazy_static! {
4      pub static ref CONSTRAINTS: RwLock<Solver> =
5          RwLock::new(Solver::new());
6  }
7
8  pub struct Solver {
9      pub constraints: HashMap<i64, Vec<ConstraintSerialized>>,
10  }
11
12  impl Solver {
13      fn new() -> Self {
14          Self {
15              constraints: HashMap::new(),
16          }
17      }
18
19      pub fn add_constraint(
20          &mut self,
21          id: i64,
22          left_operand_kind: i8,
23          left_operand_val: i64,
24          right_operand_kind: i8,
25          right_operand_val: i64,
26          predicate: i8,
27      ) {
28          self.constraints
29              .entry(id)
30              .or_insert(vec![])
31              .push(ConstraintSerialized::new(
32                  id,
33                  left_operand_kind,
34                  left_operand_val,
35                  right_operand_kind,
36                  right_operand_val,
37                  predicate,
```

```

37         ));
38     }
39     ...
40     ...
41 }

```

The next exposed runtime function call is `__symbolic_make_prepare()`, which adds a symbolic variable identifier with a tuple `<pointer address, identifier>`.

```

1 // symbolic_runtime/src/runtime.rs Rust
2
3 #[no_mangle]
4 pub extern "C" fn __symbolic_make_prepare(ptr: *mut libc::c_void,
5     addr: i64) {
6     ADDRS.write().unwrap().insert(ptr as i32, addr);
7 }

```

The next exposed runtime function call is `__make_symbolic()`, which supplies solutions from the solver.

```

1 // symbolic_runtime/src/runtime.rs Rust
2
3 #[no_mangle]
4 pub extern "C" fn __make_symbolic(typ_size: usize, ptr: *mut
5     libc::c_void) {
6     if !ptr.is_null() {
7         let id = select_id();
8         if id.is_none() {
9             return;
10        }
11        let addr = {
12            let addrs = ADDRS.read().unwrap();
13            *addrs.get(&(ptr as i32)).unwrap()
14        };
15        let id = id.unwrap();
16        if let Some(solutions) =
17            CONSTRAINTS.read().unwrap().solve(id) {
18            for (sym_addr, solution) in solutions {
19                if sym_addr == addr {
20                    match typ_size {
21                        1 => unsafe {
22                            std::ptr::write(ptr as *mut
23                                libc::c_char, solution as i8);
24                        }
25                    }
26                }
27            }
28        }
29    }
30 }

```

```

23         4 => unsafe {
24             std::ptr::write(ptr as *mut libc::c_int,
25                             solution as i32);
26         },
27         _ => unsafe {
28             std::ptr::write(ptr as *mut
29                             libc::c_long, solution);
30         },
31     }
32 }
33 }
34 }

```

It first randomly selects a state ID and retrieves the symbolic variable identifier associated with the given pointer address. This ensures the solution is supplied to the exact symbolic variable.

If a solution is found for the given constraint set, it is written into the symbolic variable's pointer. Now, let's examine how to interact with the Z3 solver to find a solution based on these constraints.

```

1  // symbolic_runtime/src/symbolic.rs Rust
2
3  fn get_symbolic_val<'ctx>(<
4      ctx: &'ctx Context,
5      var_names: &mut HashMap<i64, String>,
6      var_cnt: &mut u64,
7      kind: i8,
8      val: i64,
9  ) -> Int<'ctx> {
10      if kind == VAR_KIND {
11          if var_names.get(&val).is_none() {
12              let var_name = format!("var{}", var_cnt);
13              *var_cnt += 1;
14              var_names.insert(val, var_name);
15          }
16          let var_name = var_names.get(&val).unwrap().clone();
17          Int::new_const(&ctx, var_name.clone())
18      } else {
19          Int::from_i64(&ctx, val)
20      }
21  }
22

```

```

23 fn set_solver_timeout(ctx: &Context, solver: &Z3Solver, ms: u32) {
24     let mut params = Params::new(&ctx);
25     params.set_u32("timeout", ms);
26     solver.set_params(&params);
27 }
28
29 pub fn select_id() -> Option<i64> {
30     let constraints = CONSTRAINTS.read().unwrap();
31     let ids: Vec<_> = constraints.constraints.keys().collect();
32     if !ids.is_empty() {
33         Some(*ids.choose(&mut rand::rng()).unwrap())
34     } else {
35         None
36     }
37 }
38
39 pub fn solve(&self, id: i64) -> Option<Vec<(i64, i64)>> {
40     let cfg = Config::new();
41     let ctx = Context::new(&cfg);
42     let solver = Z3Solver::new(&ctx);
43     let mut stmts = vec![];
44     let mut syms = HashMap::new();
45     let mut var_names = HashMap::new();
46     let mut var_cnt = 0;
47     let constraints = self.constraints.get(&id).unwrap();
48     for constraint in constraints {
49         let (left_sym_val, right_sym_val) = (
50             get_symbolic_val(
51                 &ctx,
52                 &mut var_names,
53                 &mut var_cnt,
54                 constraint.left_operand_kind,
55                 constraint.left_operand_val,
56             ),
57             get_symbolic_val(
58                 &ctx,
59                 &mut var_names,
60                 &mut var_cnt,
61                 constraint.right_operand_kind,
62                 constraint.right_operand_val,
63             ),

```

```

64     );
65     let stmt;
66     match constraint.predicate {
67         PREDICATE_EQ => {
68             stmt = left_sym_val._eq(&right_sym_val);
69         }
70         PREDICATE_NE => {
71             stmt = left_sym_val._eq(&right_sym_val).not();
72         }
73         PREDICATE_SLT => {
74             stmt = left_sym_val.lt(&right_sym_val);
75         }
76         PREDICATE_SLE => {
77             stmt = left_sym_val.le(&right_sym_val);
78         }
79         PREDICATE_SGT => {
80             stmt = left_sym_val.gt(&right_sym_val);
81         }
82         PREDICATE_SGE => {
83             stmt = left_sym_val.ge(&right_sym_val);
84         }
85         _ => unreachable!("unexpected predicate: {}",
86             constraint.predicate),
87     };
88     syms.insert(constraint.left_operand_val, left_sym_val);
89     syms.insert(constraint.right_operand_val, right_sym_val);
90     stmts.push(stmt.clone());
91     solver.assert(&stmt);
92 }
93 let combined = Bool::and(&ctx,
94     &stmts.iter().collect::<Vec<_>>());
95 set_solver_timeout(&ctx, &solver, 5000);
96 match solver.check() {
97     SatResult::Sat => {
98         let model = solver.get_model().unwrap();
99         let mut solutions = vec![];
100         for (val, sym) in &syms {
101             let solution =
102                 model.eval(sym).unwrap().as_i64().unwrap();
103             if var_names.get(val).is_some() {
104                 solutions.push((*val, solution));
105             }
106         }

```

```

103      }
104      Some(solutions)
105    }
106    SatResult::Unsat => {
107      println!("UNSAT: {:?}", combined);
108      None
109    }
110    SatResult::Unknown => {
111      println!("UNKNOWN: {:?}", combined);
112      None
113    }
114  }
115 }

```

The `get_symbolic_val()` function is a utility that generates temporary symbolic variable names for use with the solver.

`set_solver_timeout()` adds a time constraint to the solver; in our implementation, this timeout is hardcoded to 5 seconds.

`select_id()` simply chooses a random state ID.

The `solve()` function contains the core logic for building statements and querying the solver.

Before we deep dive into the `solve()` body, let's first explore how to use Z3.

```

1  // example-code/symbolic/example.smt
2
3  ; Declare integer variables x and y.
4  (declare-const x Int)
5  (declare-const y Int)
6
7  ; Add constraints (assertions) to the solver.
8  ; x must be greater than or equal to 0.
9  (assert (>= x 0))
10 ; y must be greater than or equal to 0.
11 (assert (>= y 0))
12 ; The sum of x and y must be 10.
13 (assert (= (+ x y) 10))
14 ; x must be less than y.
15 (assert (< x y))
16
17 ; Check if the current set of assertions is satisfiable.
18 (check-sat)
19
20 ; If satisfiable, retrieve and print the model (variable
    assignments).

```

SMT



## 21 (get-model)

This example code illustrates how to construct a statement and query the solver.

- At lines 4-5, we declare two symbolic variables, named “x” and “y”.
- From lines 7-15, various constraints are added.
- At line 18, we check if a solution exists for the given constraints.
- Finally, at line 20, if the constraints are satisfiable, we retrieve a concrete solution.

Here is the result:

```
1 > z3 example.smt
2 sat
3 (
4   (define-fun y () Int
5     6)
6   (define-fun x () Int
7     4)
8 )
```

The solve() function encodes constraints using the Rust Z3 binding.

First, from lines 50 to 64, we create symbolic variable names. Then, from lines 66 to 86, the left and right symbolic variables are constructed into statements based on the given predicate. At line 90, each built statement is recorded. This process repeats for every constraint.

At line 92, all these statements are joined using the AND (&&) operator. Line 94 checks if the combined statement is satisfiable. If it is, the solutions for each symbolic variable are collected. Finally, these solutions are written back into the program via the `__make_symbolic()` function.

We have implemented a very basic symbolic execution engine, and there is significant room for improvement to make it more practical.

### 6.4.2.1 Homework

**Exercise 6.1** We suggest the following topics for further study and exploration:

**Topic #1: Extending Symbolic Variable Type:** Our current implementation supports symbolic variable types only for integers and characters. Readers may extend this to broader types, such as arrays.

**Topic #2: Concolic Execution:** Our current implementation is effective only for simple examples, like the working example, where all states can be identified easily without any loss. In the following example, we will introduce a genuinely difficult problem and explain why our current implementation cannot solve it.

```
1 #define H 7
2 #define W 11
3 char maze[H][W] = { "+-+---+---+",
4                      "|S|      |#|",
5                      "| |  --+ | |",
```

```

6          " | | | | |",
7          " | +-- | | |",
8          " |   |   |",
9          "+-----+---+" };
10 char* arr[28];
11 int x = 1; int y = 1; int i = 0;
12 #define ITERS 28
13 __make_symbolic(ITERS, arr);
14 while (i < ITERS) {
15     switch (arr[i]) {
16         case 'w': y--; break;
17         case 's': y++; break;
18         case 'a': x--; break;
19         case 'd': x++; break;
20         default: printf("Wrong command\n"); exit(-1);
21     }
22     if (maze[y][x] == '#') {
23         printf ("You win!\n"); exit (1);
24     }
25     if (maze[y][x] != ' ' && !((y == 2 && maze[y][x] == '|' && x > 0
&& x < W))) {
26         x = ox;
27         y = oy;
28     }
29     if (ox==x && oy==y){
30         printf("You lose\n");
31         exit(-2);
32     }
33 }

```

The example above is a maze program, abstracted from the [27] problem.

The `arr` variable holds a sequence of arrow commands. Multiple solutions exist, such as `ssssdddddwwaawdddddssssddwww`. Our current symbolic execution can identify four possible states within the switch statement. However, the core challenge lies at line 22. While we can generate 'w', 's', 'a', or 'd' for each state, the condition `maze[y][x] == '#'` is not directly tied to these four explicit states. This means the program's implicit states (like the maze layout) need to be explored by generating various sequences of arrow commands.

Consider whether this could be solved with random generation. Our current implementation can generate the four appropriate arrows ('w', 's', 'a', 'd'), but it doesn't know which arrow should be placed at which element within the sequence. If we were to generate these randomly, the probability of reaching line 23 would be roughly  $0.25^{28} = 1.3877787807814457e-17$  —an utterly impractical solution in the real world.

Here, the concept of *state explosion* becomes apparent. Randomly mixing states often fails to satisfy complex branch conditions. To mitigate state explosion, *state pruning* can be adopted, where useless states are discarded. For example, a command like “sssssss” will not be a maze solution, and the program will likely hit line 30 at the fifth ‘s’. At that point, having covered the line 30 branch, we can prioritize other uncovered branches. This suggests using backtracking to quickly revert to a state before a previous coverage point was found. We would then go back to the command “ssss” and try generating ‘w’, ‘a’, or ‘d’ for the fifth character. This would yield three new commands: “sssw”, “sssa”, or “sssd”. If “sssw” and “sssd” hit line 30 again, we would prune them from the set of states and continuously proceed with the remaining viable commands, repeating this process.

For this homework, it’s sufficient to understand this problem. Passionate readers may attempt to implement a symbolic execution engine capable of solving the maze problem within a reasonable time budget. Implementing state pruning and scheduling techniques would be necessary to reduce useless exploration and save time. Some readers might realize at this point that the basis of state pruning often relies on coverage. If a state has been explored sufficiently, and it’s confirmed that it won’t lead to further new coverage, that state can be safely removed to explore uncovered code quickly.



## Part VI

# Delta Debugging

	<b>Introduction to Delta Debugging.</b>	
<b>7</b>	<b>159</b>	
7.1	Introduction . . . . .	159
7.2	Simple Idea (Binary Search) . . . . .	159
7.3	Local Minmum and Global Minimum .	160
7.4	1-minimal Algorithm . . . . .	161
7.5	Implementation . . . . .	161
7.6	Testing . . . . .	164
7.7	Fuzzer Integration . . . . .	166



# “Buggy input found. Is it the minimized one?”

Input1: '0x1aa2bbac1addabb8b1b2'	✗
Input2: '0x1aa2bbac1add'	✗
Input3: '0x1aa2'	✓
Input4: '0x1aa2bba'	✗

## 7. Introduction to Delta Debugging

### 7.1 Introduction

Let's assume our fuzzer has found a crash with an input 100 bytes long. A programmer can reproduce and confirm this bug using that specific seed. However, this isn't the final step. While we've identified a crash-inducing input, is it truly minimal? There might be another input that causes the same crash but is significantly shorter. A minimized input greatly aids the programmer in understanding why this input led to a bug. Simply put, when comparing a 3-byte bug-inducing input with a 100-byte one, which is easier to debug? Absolutely the former.

This brings us to Delta Debugging (DD) [28], a methodology for finding minimal failure-inducing inputs. The ultimate goal of this chapter is to integrate DD into our fuzzer, enabling it to automatically reduce crash-inducing seeds.

### 7.2 Simple Idea (Binary Search)

Imagine our target input contains the characters '1' and '7'. A crashing input might be "12345678" because it includes both.

Perhaps we can apply a binary search approach to find the minimal failure-inducing input. Let's explore that.

Step	Test	Result
0	1 2 3 4 5 6 7 8	
1	1 2 3 4 . . . .	Pass
2	. . . . 5 6 7 8	Pass
3	1 2 . . . . .	Pass
4	. . 3 4 . . . .	Pass
5	. . . . 5 6 . .	Pass
6	. . . . . 7 8	Pass
7	1 . . . . . .	Pass
8	. 2 . . . . .	Pass

Step	Test	Result
9	.. 3 .....	Pass
10	... 4 .....	Pass
11	.... 5 ...	Pass
12	..... 6 ..	Pass
13	..... 7 .	Pass
14	..... 8	Pass

Even with a failure-inducing input, binary search failed to find the bug. This is because the process doesn't account for the **“complement”**. The complement is the combined set of elements excluding the current delta. For instance, in Step 3, the complement of the input (“12”) is “345678”. Thus, formally, the complement is defined as:

$$c_x = \text{failure-inducing input (e.g., 12345678 in the example)} \quad (7.1)$$

$$\Delta_i = \text{delta of } c_x \quad (7.2)$$

$$\nabla_i(\text{Complement}) = c_x - \Delta_i \quad (7.3)$$

By definition, in Step 3, the complement is “12345678” - “12” = “345678”. Therefore, considering the complement is key to DD, and we will now introduce the **1-minimal algorithm**.

### 7.3 Local Minmum and Global Minimum

Finding a global minimum is a time-consuming and challenging problem. In the real world, finding a local minimum is usually sufficient.

- **n-minimal**: A failure-inducing input where the bug disappears if any n elements are removed.
- **1-minimal**: A failure-inducing input where the bug disappears if any single element is removed.

Finding an n-minimal input is more challenging than finding a 1-minimal input because the 1-minimal algorithm is greedy; it stops as soon as the bug disappears.

In the previous example, with the input “12345678”, both n-minimal and 1-minimal algorithms successfully find the global minimum, which is “17”. Now, let's explore a more complex buggy condition.

Assume the buggy condition is as follows:

- If the input contains both ‘A’ and ‘B’, a bug is triggered.
- If the input contains both ‘A’ and ‘C’, a bug is triggered.

Given the input “ABCDE”, let's assume “ABC” is an intermediate result of the 1-minimal algorithm. In the next step, when it attempts to delete ‘A’, the resulting “BC” does not report a bug because it satisfies neither condition 1 nor condition 2. If either ‘B’ or ‘C’ were removed, the results (‘AC’ or ‘AB’) would trigger a bug. However, since the algorithm found a non-triggering condition (“BC”), the 1-minimal algorithm will incorrectly conclude that the minimized result is “ABC”. The correct answers, however, are “AB” or “AC”. Therefore, the 1-minimal algorithm does not guarantee a global minimum.



## 7.4 1-minimal Algorithm

We'll now introduce the 1-minimal algorithm from the paper [29].

The high-level idea is straightforward: the algorithm repeatedly divides the input into a **delta (subset)** and its **complement**, then runs a test (or oracle) to see if a bug appears. This process continues until no bug is found or the granularity can't be refined further.

We attach the intuitive algorithm here from the paper, so called, "**ddmin**".

### Minimizing Delta Debugging Algorithm

Let  $test$  and  $c_{\mathbf{x}}$  be given such that  $test(\emptyset) = \checkmark \wedge test(c_{\mathbf{x}}) = \mathbf{X}$  hold.

The goal is to find  $c'_{\mathbf{x}} = ddmin(c_{\mathbf{x}})$  such that  $c'_{\mathbf{x}} \subseteq c_{\mathbf{x}}$ ,  $test(c'_{\mathbf{x}}) = \mathbf{X}$ , and  $c'_{\mathbf{x}}$  is 1-minimal.

The minimizing Delta Debugging algorithm  $ddmin(c)$  is

$ddmin(c_{\mathbf{x}}) = ddmin_2(c_{\mathbf{x}}, 2)$  where

$$ddmin_2(c'_{\mathbf{x}}, n) = \begin{cases} ddmin_2(\Delta_i, 2) & \text{if } \exists i \in \{1, \dots, n\} \cdot test(\Delta_i) = \mathbf{X} \text{ ("reduce to subset")} \\ ddmin_2(\nabla_i, \max(n-1, 2)) & \text{else if } \exists i \in \{1, \dots, n\} \cdot test(\nabla_i) = \mathbf{X} \text{ ("reduce to complement")} \\ ddmin_2(c'_{\mathbf{x}}, \min(|c'_{\mathbf{x}}|, 2n)) & \text{else if } n < |c'_{\mathbf{x}}| \text{ ("increase granularity")} \\ c'_{\mathbf{x}} & \text{otherwise ("done").} \end{cases}$$

where  $\nabla_i = c'_{\mathbf{x}} - \Delta_i$ ,  $c'_{\mathbf{x}} = \Delta_1 \cup \Delta_2 \cup \dots \cup \Delta_n$ , all  $\Delta_i$  are pairwise disjoint, and  $\forall \Delta_i \cdot |\Delta_i| \approx |c'_{\mathbf{x}}|/n$  holds.

The recursion invariant (and thus precondition) for  $ddmin_2$  is  $test(c'_{\mathbf{x}}) = \mathbf{X} \wedge n \leq |c'_{\mathbf{x}}|$ .

Figure 7.1: 1-Minimal Algorithm

This algorithm is quite intuitive. Let's break it down step by step.

The  $ddmin()$  function kicks off the process with a fixed granularity parameter of '2'. This '2' dictates the initial granularity, meaning the algorithm first divides the input into segments that are half the size of the original input.

The  $ddmin_2()$  function's initial task is to split the input into delta and complement sets. After this, we iterate through all deltas and complements. We'll start by executing the delta set.

If a failure is detected while testing with a delta set, the function recursively calls itself with the current delta and a granularity of 2 (the same as the initial case).

If a failure is found during testing with the complement set, it recursively calls itself with the current complement set and  $\max(n-1, 2)$ . This means the granularity is slightly adjusted if  $n$  is greater than 2, but the minimum granularity remains 2, mirroring the first case.

Finally, if no error is found in either the delta or complement set, the algorithm attempts to increase the granularity to explore more test cases.

The algorithm stops if the function flow doesn't fall into any of these three conditions.

## 7.5 Implementation

We start implement  $ddmin()$  as a library first and test. Once we understood by implementing the algorithm with library first, and we will port it into the fuzzer.

```
1 // delta_debugging/src/lib.rs
```

 Rust

```
2  pub type Data = Vec<u8>;
3
4  #[derive(PartialEq, Debug)]
5  pub enum TestResult {
6      Pass,
7      Fail,
8  }
9
10 pub type TestFn = Box<dyn Fn(&Data) -> TestResult>;
11
12 pub fn ddmin(data: &Data, test: TestFn) -> Data {
13     do_ddmin(data, 2, test)
14 }
15
16 fn do_ddmin(data: &Data, n: usize, test: TestFn) -> Data {
17     let (delta_set, complement_set) = split(data, n);
18     for delta in &delta_set {
19         if test(delta) == TestResult::Fail {
20             if delta.len() == 1 {
21                 return delta.to_vec();
22             }
23             return do_ddmin(delta, 2, test);
24         }
25     }
26     for complement in &complement_set {
27         if test(complement) == TestResult::Fail {
28             return do_ddmin(complement, max(n - 1, 2), test);
29         }
30     }
31     if n < data.len() {
32         return do_ddmin(data, min(data.len(), 2 * n), test);
33     }
34     data.to_vec()
35 }
36
37 /// Return delta and complement
38 pub fn split(data: &Data, n: usize) -> (Vec<Data>, Vec<Data>) {
39     if n == 0 {
40         return (Vec::new(), Vec::new());
41     }
42     let data_len = data.len();
```

```
43     let exact_chunk_size = data_len / n;
44     let remainder = data_len % n;
45
46     let mut delta_boundaries = Vec::new();
47     let mut cur_pos = 0;
48
49     for i in 0..n {
50         let mut chunk_size = exact_chunk_size;
51         if i < remainder {
52             chunk_size += 1;
53         }
54         if cur_pos + chunk_size > data_len {
55             break;
56         }
57         delta_boundaries.push((cur_pos, cur_pos + chunk_size));
58         cur_pos += chunk_size;
59     }
60
61     let delta_set: Vec<Data> = delta_boundaries
62         .iter()
63         .map(|&(start, end)| data[start..end].to_vec())
64         .collect();
65
66     let mut complement_set: Vec<Data> = Vec::new();
67     for i in 0..delta_set.len() {
68         let (start, end) = delta_boundaries[i];
69         let mut complement_bytes_raw = Vec::new();
70         if start > 0 {
71             complement_bytes_raw.extend_from_slice(&data[0..start]);
72         }
73         if end < data_len {
74             complement_bytes_raw.extend_from_slice(
75                 &data[end..data_len]);
76         }
77         complement_set.push(complement_bytes_raw);
78     }
79     let filtered_complements: Vec<Data> = complement_set
80         .into_iter()
81         .filter(|c| !delta_set.contains(c))
82         .collect();
83     (delta_set, filtered_complements)
```

```
84 }
```

The implementation closely mirrors the algorithm. The primary recursive call function is defined on line 16. Line 12 defines `ddmin2()`, which serves as the algorithm's entry function. It's worth noting that the `test` function acts as the real-world oracle, determining whether an input is failure-inducing. Therefore, we treat `test` as a first-class function that must be defined by the caller. The `split()` function returns both the delta and complement sets.

## 7.6 Testing

The following functions are for Delta Debugging (DD) testing, with multiple test cases. The first case uses the input demonstrated in the paper's working example.

```
1 // delta_debugging/tests/dd_test.rs Rust
2
3 #[test]
4 fn test_ddmin() {
5     let tcs = [
6         (
7             String::from("12345678"),
8             String::from("178"),
9             String::from("178"),
10        ),
11        (
12            String::from("12345678"),
13            String::from("1<78"),
14            String::from("12345678"),
15        ),
16        (
17            String::from("int a = 1; int b = 2; assert(a == b);"),
18            String::from("assert(a == b);"),
19            String::from("assert(a == b);"),
20        ),
21        (
22            String::from("This is a test string with a problematic $
23                           character inside."),
24            String::from("$"),
25            String::from("$"),
26        ),
27        (
28            String::from("function calculate_sum(a, b) { return a +
29                           b; print(\"calc\" ) }"),
30            String::from("print(\"calc\"),
31            String::from("print(\"calc\"),
```

```

30     ),
31     (
32         String::from(
33             "{ 'data': [1, 2, 3, 'value', {'key': 'nested'},
34             'extra', {'bug': 'missing_brace'}] }",
35         ),
36         String::from("missing_brace"),
37         String::from("'missing_brace'"),
38     );
39     for tc in tcs {
40         let (input, fail_inducing_input, expected) = (tc.0, tc.1,
41         tc.2);
42         test(input, fail_inducing_input, expected.into_bytes());
43     }
44 }
45 fn test(input: String, fail_inducing_input: String, expected: Data)
46 {
47     let oracle = make_oracle(fail_inducing_input.into_bytes());
48     let minimized = ddmin(&input.into_bytes(), oracle);
49     println!(
50         " {:?} == {:?}",
51         byte_to_str(&minimized),
52         byte_to_str(&expected)
53     );
54     assert_eq!(minimized, expected);
55 }
56 fn byte_to_str(data: &Data) -> String {
57     String::from_utf8(data.to_vec()).unwrap()
58 }
59
60 fn make_oracle(fail_inducing_input: Data) -> Box<dyn Fn(&Data) ->
61     TestResult> {
62     Box::new(move |data: &Data| {
63         let mut fail_inducing_cnt = HashMap::new();
64         for v in &fail_inducing_input {
65             *fail_inducing_cnt.entry(v).or_insert(0) += 1;
66         }
67         let mut input_cnt = HashMap::new();
68         for v in data {

```

```

68         *input_cnt.entry(v).or_insert(0) += 1;
69     }
70     for (key, &needed) in &fail_inducing_cnt {
71         let available =
72             input_cnt.get(key).copied().unwrap_or(0);
73         if available < needed {
74             println!("pass: {:?}", byte_to_str(data));
75             return TestResult::Pass;
76         }
77         println!("failed: {:?}", byte_to_str(data));
78         TestResult::Fail
79     })
80 }

```

The test input comprises: the original long failure-inducing input, its exact failure-inducing answer, and the expected 1-minimal output. We've defined an oracle that returns `TestResult::Pass` if the given input contains the failure-inducing elements, and `TestResult::Fail` otherwise.

If readers wish to observe the step-by-step progress, they can run the test with the command: `cargo test test_ddmin -- --exact --nocapture`.

## 7.7 Fuzzer Integration

Having studied and implemented the 1-minimal algorithm, it's time to integrate it into our fuzzer. Let's first outline the necessary steps:

- **Oracle Definition:** Our fuzzer's campaign already serves as the oracle. The oracle, in this context, is the fuzzing process itself, testing the input that's being minimized by the 1-minimal procedure.
- **Seed Minimization:** When a crash-inducing seed is found, we enter a minimization loop—the 1-minimal procedure. This pauses the generation of new seeds. Once the minimized seed is identified, we store it in the crash directory and then resume the original fuzzing process.

```

1  // fuzzer/src/fuzzer.rs Rust
2
3  fn is_crash(&self, status: i32) -> bool {
4      status != PROCESS_EXIT_NORMAL as i32 && status != SIGKILL
5  }
6
7  fn oracle(
8      &mut self,
9      child_stdin: &mut Option<ChildStdin>,
10     timeout: Duration,
11     seed: &Seed,

```

```
12 ) -> Result<TestResult> {
13     self.feed_seed(child_stdin, seed)?;
14     self.wakeup_forkserver(HostSend::Wakeup(timeout.as_secs()));
15     let status = self.wait_forkserver();
16     self.clear_new_coverage();
17     self.clear_visited_edges();
18     if self.is_crash(status) {
19         Ok(TestResult::Fail)
20     } else {
21         Ok(TestResult::Pass)
22     }
23 }
24
25 fn ddmin(
26     &mut self,
27     child_stdin: &mut Option<ChildStdin>,
28     timeout: Duration,
29     seed: &Seed,
30 ) -> Result<Seed> {
31     self.do_ddmin(child_stdin, timeout, seed, 2)
32 }
33
34 fn do_ddmin(
35     &mut self,
36     child_stdin: &mut Option<ChildStdin>,
37     timeout: Duration,
38     seed: &Seed,
39     n: usize,
40 ) -> Result<Seed> {
41     let (delta_set, complement_set) =
42         split(&seed.get_input().to_vec(), n);
43     for delta in &delta_set {
44         let delta_seed = Seed::new(delta.to_vec(), 0);
45         if let Ok(test_result) = self.oracle(child_stdin, timeout,
46             &delta_seed) {
47             if test_result == TestResult::Fail {
48                 if delta.len() == 1 {
49                     return Ok(delta_seed);
50                 }
51             }
52         }
53     }
54     return self.do_ddmin(child_stdin, timeout,
55         &delta_seed, 2);
56 }
```

```

51     }
52 }
53 for complement in &complement_set {
54     let complement_seed = Seed::new(complement.to_vec(), 0);
55     if let Ok(test_result) = self.oracle(child_stdin, timeout,
56         &complement_seed) {
57         if test_result == TestResult::Fail {
58             return self.do_ddmin(child_stdin, timeout,
59                 &complement_seed, max(n - 1, 2));
60         }
61     }
62     let seed_len = seed.get_input().len();
63     if n < seed_len {
64         return self.do_ddmin(child_stdin, timeout, seed,
65             min(seed_len, 2 * n));
66     }
67     Ok(seed.clone())
68 }

```

```

1  // fuzzer/src/fuzzer.rs Rust
2
3  pub fn run(&mut self, program_file: &str) -> Result<FuzzResult> {
4      ...
5      ...
6      // crash found
7      if self.is_crash(status) {
8          if let Ok(minimized) = self.ddmin(&mut child_stdin, timeout,
9              &seed) {
10             self.crashes.insert(minimized.clone());
11             minimized.to_file(self.crashes.len());
12             self.send(FuzzShot::Crash(CrashInfo::new(
13                 self.crashes.len(),
14                 seed,
15                 minimized,
16             )));
17         }
18     }
19 }

```

The `oracle()` function utilizes the fuzzing procedure as an oracle. The oracle result is then passed to the 1-minimal procedure for further minimization.



### 7.7.1 Homework

**Exercise 7.1** We suggest the following topics for further study and exploration:

- **Optimization:** Readers can explore state-of-the-art papers for more optimized solutions to find minimized failure-inducing inputs.



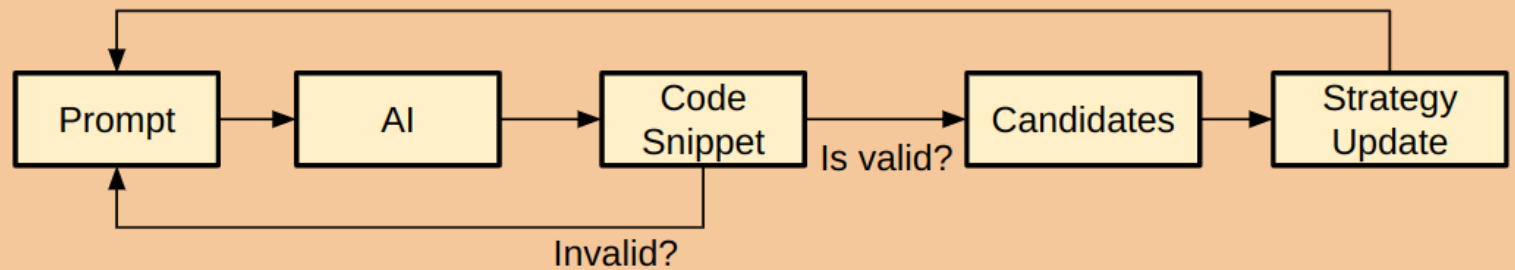
## Part VII

# LLM-based Synthesis

<b>8</b>	<b>LLM-based Synthesis . . . . .</b>	<b>173</b>
8.1	Hello, AI . . . . .	173
8.2	Introduction to LLM-based Fuzzer . . .	174
8.3	Introduction to Program Synthesis . . .	176
8.4	LLM-based Synthesis . . . . .	177
8.5	Retrospect . . . . .	184



# “Leverage an LLM to build a program that generates programs”



## 8. LLM-based Synthesis

### 8.1 Hello, AI

The ubiquity of AI is no longer a major hurdle. Across almost every field, interactive Large Language Model (LLM [30]) are widely used and provide excellent assistance for daily tasks.

My first real experience with AI was in 2019. At that time in South Korea, there were very few AI lectures offered in university courses. Computer vision was the main topic for AI integration, and Convolutional Neural Network (CNN [31]) were a hot topic, with the MNIST [32] classification problem often serving as a starting point. Simultaneously, in the academic world, Natural Language Processing (NLP [33]) was a hot topic, with models like Long Short-Term Memory (LSTM [34]) being a primary example for tasks like language translation.

Back then, I didn't expect AI to become so dominant, especially not how strongly it would become tied to our daily lives.

Nowadays, when a new keyword or concept comes to my attention, I prefer to ask an LLM rather than searching on Google. When I used to Google, I had to expend a lot of effort to filter out correct from incorrect information, identify the key points, understand the motivation, and grasp the output. Since I lacked the foundational knowledge (which was why I was searching in the first place), this process was very time-consuming. When I first started asking LLMs, I was surprised that the answers to my questions seemed correct and directly addressed the major concepts I was curious about. Moreover, the interactive nature of LLMs allows me to ask follow-up questions. It's fascinating to reflect on how AI has integrated into daily life, much like the experience I've shared.

As a senior CS student in 2019, I didn't anticipate this success. Today, I sometimes wonder what the next hot, winning, and game-changing technology will be. Personally, I believe the AI model itself is not the only important part. I now look at the entire AI ecosystem to see which parts can be improved, where I can contribute, which parts are connected to computer science (software), and what seems interesting and fun.

My current and graduate majors were not in AI. However, learning about other fields has led to new discoveries in my personal research, provided a new perspective, and reignited the joy of programming. So, why not take some time today to explore an AI-related product, technology, or trend? It can also be beneficial to research other fields,

even if the keyword isn't "AI," and see a different area from your daily work. It may give you a perspective you couldn't have imagined before.

In this chapter, our goal is to implement an **LLM-based synthesis tool** to generate test code for our implementation.

## 8.2 Introduction to LLM-based Fuzzer

Academic research is actively exploring how to integrate LLMs. For example, a referenced paper on an LLM-based fuzzer, [35], successfully found new, previously undiscovered bugs. In my opinion, this type of fuzzer is particularly effective at finding bugs in new features. A prime use case would be testing new syntax or libraries of a programming language, new API usages, and so on. This fuzzer typically consists of two main components.

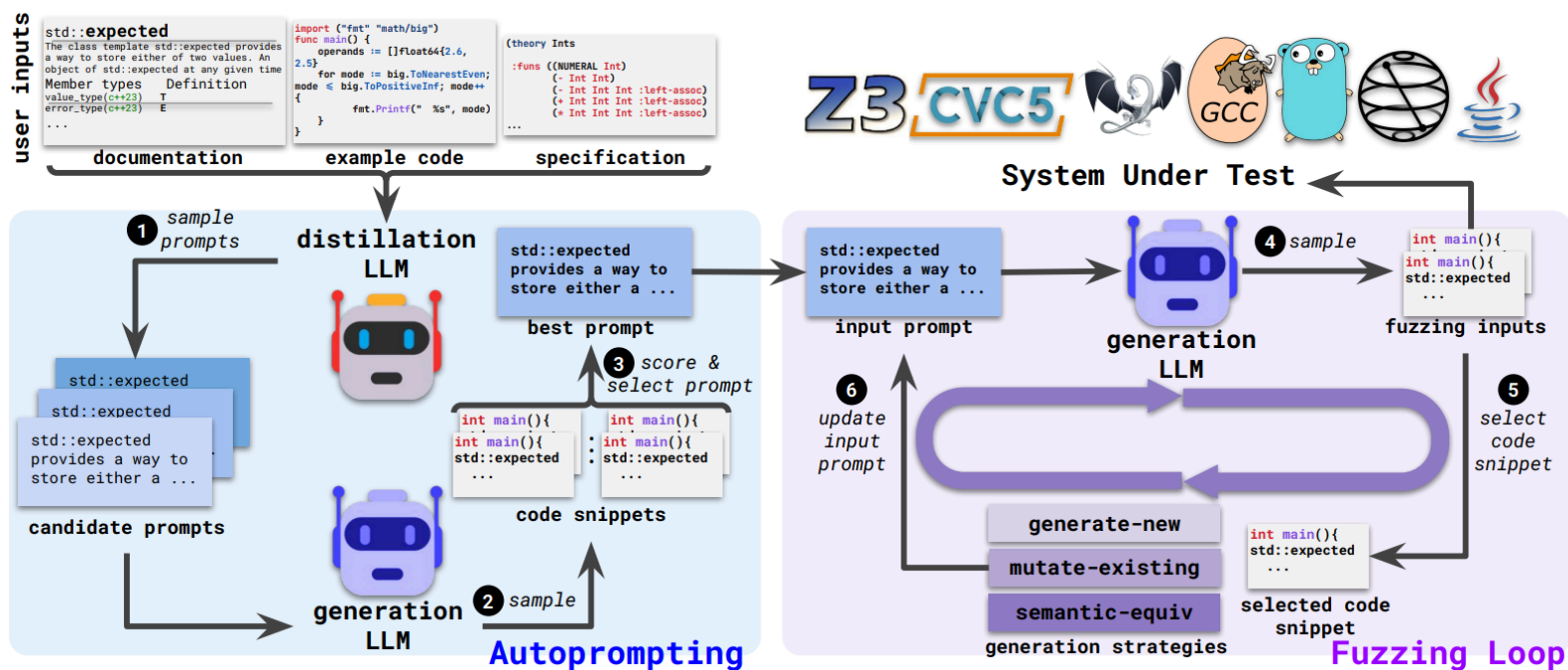


Figure 8.1: Overview of Fuzz4All

The architecture, which is based on a figure from the paper, has two major components: **Autoprompting** and the **Fuzzing Loop**.

### 8.2.1 Autoprompting

Autoprompting is a one-time initialization task that runs before the fuzzing loop begins. It takes documentation, example code, and specifications as input. These inputs are fed to an LLM, which is guided to generate a concise, abstracted set of statements. A scoring function then assigns scores to these statements based on various metrics, such as code coverage or bug-finding potential. The paper's approach specifically assigns a higher score if a generated code snippet is accepted and runnable by the target System Under Test (SUT). This process ultimately yields a set of high-quality candidate prompts to be used in the fuzzing loop.

---

**Algorithm 1:** Autoprompting for fuzzing
 

---

```

1 Function Autoprompting:
  Input : userInput, numSamples
  Output: inputPrompt

2   greedyPrompt  $\leftarrow \mathcal{M}_{\mathcal{D}}$  (userInput, APInstruction, temp=0)
3   candidatePrompts  $\leftarrow$  [ greedyPrompt ]
4   while |candidatePrompts| < numSamples do
5     prompt  $\leftarrow \mathcal{M}_{\mathcal{D}}$  (userInput, APInstruction, temp=1)
6     candidatePrompts  $\leftarrow$  candidatePrompts + [ prompt ]
7   inputPrompt  $\leftarrow \arg \max_{p \in \text{candidatePrompts}} \text{Scoring}(\mathcal{M}_{\mathcal{G}}(p), \text{SUT})$ 
8   return inputPrompt

```

---

Figure 8.2: AutoPrompting

The autoprompting algorithm, which is taken from the paper, is described above.

It begins by generating an initial “greedy” prompt from a distillation model with a temperature of zero. This setting makes the output highly conservative. The process continues until a sufficient number of candidate prompts have been collected.

To explore a wider range of possibilities, the algorithm then generates additional prompts using the same format but with the temperature set to one. This encourages more creative and diverse outputs. Finally, the best prompt is selected using the `argmax()` function, which finds the highest-scoring candidate based on the defined scoring function.

### 8.2.2 Fuzzing Loop

The prompt generated during the Autoprompting phase is fed into a generation LLM, which yields a code snippet based on the input. This generated code snippet is then tested against the SUT. Mutation is a key part of this process, and there are three mutation strategies that we will detail shortly. This entire loop is repeated until the given time budget is exhausted.

**Algorithm 2: Fuzzing loop**

```

1 Function FuzzingLoop:
  Input :inputPrompt, timeBudget
  Output: bugs
2   genStrats  $\leftarrow$  [ generate-new, mutate-existing,
    semantic-equiv ]
3   fuzzingInputs  $\leftarrow$   $\mathcal{M}_{\mathcal{G}}$  (inputPrompt + generate-new)
4   bugs  $\leftarrow$  Oracle (fuzzingInputs, SUT)
5   while timeElapsed < timeBudget do
6     example  $\leftarrow$  sample (fuzzingInputs, SUT)
7     instruction  $\leftarrow$  sample (genStrats)
8     fuzzingInputs  $\leftarrow$   $\mathcal{M}_{\mathcal{G}}$  (inputPrompt + example +
    instruction)
9     bugs  $\leftarrow$  bugs + Oracle (fuzzingInputs, SUT)
10  return bugs

```

Figure 8.3: Fuzzing Loop

As described in the algorithm, there are three mutation strategies: new-generation, mutate-existing, and semantic-equivalent.

The new-generation strategy involves a prompt that creates a program from scratch. This is the strategy used to start the initial code snippet generation (Also, being used in the loop continuously). During the fuzzing loop, a strategy is selected randomly.

The mutate-existing strategy prompts the LLM to generate a mutated code snippet based on a previous generation. The semantic-equivalent strategy, on the other hand, guides the LLM to create a semantically equivalent code snippet. The generated code snippets from these strategies are fed to the SUT, and the process is repeated.

This scheme is straightforward to understand, yet its effectiveness is remarkable. It has successfully found bugs in compilers like GCC and Golang, as detailed in [36]. Readers interested in the paper’s experiments are encouraged to read the full document.

In summary, we have studied the power and motivation of LLMs, how they are integrated into fuzzing, and how they are used. Inspired by this, we will implement an LLM-based synthesis tool to generate test code snippets for our OOB detector (from Chapter 2, ASan).

### 8.3 Introduction to Program Synthesis

Program synthesis [37] is the task of automatically generating a program from a high-level specification. The synthesis tool takes a requirement and attempts to generate code that satisfies it. For example, to create a C function that adds two integers, our specification would be the function signature: `int add(int, int);`.

The main challenge is implementing the function body. To do this, the synthesizer needs a formal grammar or a set of rules it can use to construct and combine code elements. For instance:

1 a	Spec
2 b	



```

3  +
4  -
5  return
6  ;

```

Based on this grammar, the synthesizer will try to generate code snippets such as:

```

1  a b +
2  a b -
3  a return b
4  ...
5  ...

```

To minimize the generation of invalid syntax, we can guide the synthesizer with defined constraints. For example: The variables `a` and `b` must be located between an operator, such as `+` or `-`. The last statement should start with `return`.

Finally, we should provide input-output examples such as:

1 examples	Input-output
2 1 2 3	
3 2 3 5	
4 ...	

With these more specific guides, the synthesizer will eventually generate the intended function, `int add(int a, int b) { return a+b; }`. However, it might also generate semantically equivalent but more complex code, such as `int add(int a, int b) { return a+b-a-b+a+b; }`. To avoid generating unnecessary code and save time, branch pruning is needed. If you are interested in traditional program synthesis, I recommend exploring university courses.

## 8.4 LLM-based Synthesis

This chapter's main topic is the implementation of a synthesizer that generates test programs for our OOB detector. This tool can be extended to generate other programs that target specific types of bugs.

As we saw, traditional program synthesis finds it difficult to generate even a simple function like “add.” In contrast, our tool will generate more complex programs that contain two kinds of bugs: “buffer overflow” and “use-after-free”.

The synthesizer is implemented in just 114 lines of code. For this model, we've removed the Autoprompting phase and replaced it with a hand-written initial prompt, which serves similarly as an output of Autoprompting. Our initial prompt is as follows:

```

1 Requirements:
2 - Do not use any file-related functionality.
3 - Only use standard library functions: malloc, free, and strcpy.
4 - Use the function signature: void main().
5 - Add a brief comment next to the buggy line indicating the type of
  bug.

```

- ```

6 - Provide only one function, with no additional explanation or
  output.
7 - **Include only one bug type per generated code.**

```

Here are the constraints we used to guide the LLM:

- Security: For security reasons, we instructed the LLM not to use any file-related functionality or system calls.
- Standard Library: We guided the LLM to use standard library functions such as `malloc`, `free`, and `strcpy`.
- Main Function: The main function's signature should be `void main()`.
- Bug Indicators: The LLM is instructed to add a brief comment indicating the intended bug type on the line where the buggy line is instrumented.
- Conciseness: The generated code should be concise, containing only a single function with no detailed comments per line.
- Single Bug Type: Each generated code snippet must contain only one bug type and should not mix two kinds of bugs simultaneously.
- 

Here are the mutation strategy prompts (we use the same three rules, but different prompts):

- new-generation-bof: "Please create a program that triggers buffer overflow bug."
- new-generation-uaf: "Please create a program that triggers use-after-free bug."
- mutate-existing: "Please create a mutated program that modifies the previous generation."
- semantic-equiv: "Please create a semantically equivalent program to the previous generation."

Here is the complete synthesizer code:

```

1  //synthesis/asan_model.py
2
3  from transformers import (
4      AutoModelForCausalLM,
5      AutoTokenizer,
6      BitsAndBytesConfig,
7      StoppingCriteriaList,
8      StoppingCriteria,
9  )
10 import torch
11 import random
12 import re
13 import os
14 import subprocess
15
16 def extract_first_code_block(text):
17     pattern = r'\s*\s*\s*(?:\w+)?\n?(.*?)\s*\s*\s*'
18     match = re.search(pattern, text, re.DOTALL)

```

```
19     if match:
20         return match.group(1).strip()
21     return ""
22
23 def gen_model():
24     model_name = "Qwen/Qwen2.5-Coder-7B-Instruct"
25     quantization_config = BitsAndBytesConfig(
26         load_in_4bit=True,
27         bnb_4bit_compute_dtype=torch.float16,
28         bnb_4bit_quant_type="nf4",
29         bnb_4bit_use_double_quant=True,
30     )
31     tokenizer = AutoTokenizer.from_pretrained(model_name)
32     model = AutoModelForCausalLM.from_pretrained(
33         model_name,
34         quantization_config=quantization_config,
35         device_map="auto",
36         trust_remote_code=True,
37     )
38     print(f"Model loaded:
39         {torch.cuda.memory_allocated()/1024**3:.2f}GB")
40     return model, tokenizer
41
42 def gen_prompt(tokenizer, model_device, user_prompt):
43     messages = [
44         {"role": "system", "content": "You are a helpful coding
45         assistant"},
46         {"role": "user", "content": user_prompt}
47     ]
48     text = tokenizer.apply_chat_template(messages, tokenize=False,
49         add_generation_prompt=True)
50     return tokenizer([text], return_tensors="pt").to(model_device)
51
52 def execute(model, tokenizer, input_prompt):
53     with torch.no_grad():
54         generated_ids = model.generate(
55             **input_prompt,
56             max_new_tokens=256,
57             do_sample=True,
58             repetition_penalty=1.0,
59             temperature=1.0,
60             top_p=0.7,
```

```
58         pad_token_id=tokenizer.eos_token_id,
59     )
60     generated_ids = [output_ids[len(input_ids):] for input_ids,
61                     output_ids in zip(input_prompt.input_ids, generated_ids)]
62     response = tokenizer.batch_decode(generated_ids,
63                                     skip_special_tokens=True)[0]
64     return response
65
66 def is_valid_code(c_code):
67     result = subprocess.run(
68         ['clang', '-x', 'c', '-'],
69         input=c_code,
70         text=True,
71         capture_output=True
72     )
73     if result.returncode == 0:
74         os.remove("./a.out")
75     return result.returncode == 0
76
77 def gen_example(n_samples=5):
78     model, tokenizer = gen_model()
79     guide_prompt = '''
80     Requirements:
81     - Do not use any file-related functionality.
82     - Only use standard library functions: malloc, free, and
83       strcpy.
84     - Use the function signature: void main().
85     - Add a brief comment next to the buggy line indicating the
86       type of bug.
87     - Provide only one function, with no additional explanation
88       or output.
89     - **Include only one bug type per generated code.**
90     '''
91     newgen_bof_prompt = "Please create a program that triggers
92     buffer overflow bug."
93     newgen_uaf_prompt = "Please create a program that triggers use-
94     after-free bug."
95     se_prompt = "Please create a semantically equivalent program to
96     the previous generation"
97     mutate_prompt = "Please create a mutated program that modifies
98     the previous generation"
99     gen_stats = [newgen_bof_prompt, newgen_uaf_prompt,
100                 mutate_prompt, se_prompt]
```

```

92     # generated mutated fuzzing inputs
93     fuzzing_inputs = []
94     while len(fuzzing_inputs) < n_samples:
95         instruction = random.choice(gen_stats)
96         print(f"generating ...{len(fuzzing_inputs)}")
97         print(f"selected instruction: {instruction}")
98         prompt_id = gen_prompt(tokenizer, model.device,
99                                guide_prompt + "\n" + instruction)
100        fuzzing_input = execute(model, tokenizer, prompt_id)
101        fuzzing_input = extract_first_code_block(fuzzing_input)
102        if is_valid_code(fuzzing_input):
103            fuzzing_inputs.append(fuzzing_input)
104        else:
105            print(f"compile failed: \n {fuzzing_input}")
106    return fuzzing_inputs
107
108 def write_c_files(fuzzing_inputs):
109     os.makedirs("generated", exist_ok=True)
110     for i, code in enumerate(fuzzing_inputs):
111         filename = f"generated/{i + 1}.c"
112         with open(filename, 'w') as f:
113             f.write(code)
114
115 n_samples = 5
116 examples = gen_example(n_samples)
117 write_c_files(examples)

```

Note on line 17: We added an escape character ‘\’ in this code block because the character overlaps with the book’s code annotation. This character is not present in the original code.

We use the ‘Qwen/Qwen2.5-Coder-7B-Instruct’ model [38]. Its 7B parameters are a reasonable scale to be executed on a home laptop with quantization enabled.

The `gen_model()` function initializes the model and tokenizer and sets up quantization. Model quantization converts the type of the model’s weights to a lower bit, which makes it faster and more lightweight at the trade-off of a slight degradation in accuracy.

The `execute()` function invokes inference with a temperature of 1.0, limits the token size to 256, and returns the generated code snippet.

The `gen_example()` function is the main loop. It takes a number of samples, which specifies how many code snippets to generate. First, it defines the prompts for guidance and mutation. In the for loop, it randomly selects one of the mutation strategies, executes the prompt, and checks if the generated code is a valid program through the `is_valid_code()` function. Finally, the generated code is stored to disk via the `write_c_files()` function.

Let’s run the synthesizer (Author’s GPU spec is NVIDIA GeForce RTX 4060 Laptop GPU)



```
36 selected instruction: Please create a program that triggers buffer
    overflow bug.
37 generating ...3
38 selected instruction: Please create a program that triggers buffer
    overflow bug.
39 generating ...4
40 selected instruction: Please create a program that triggers use-
    after-free bug.
```

If a generated code snippet is not compilable, the output will show the code like the one above. The reason for the failure in this case is the missing declaration of `#include <stdio.h>`. The full error message is shown below:

failure.c:21:5: note: include the header `<stdio.h>` or explicitly provide a declaration for `'printf'`.

A directory named `'generated'` will be created, containing five C files: `1.c`, `2.c`, `3.c`, `4.c`, and `5.c`. Let's examine the contents of `1.c`.

```
1  #include <stdlib.h>
2  #include <string.h>
3
4  void main() {
5      char *str = (char *)malloc(20 * sizeof(char)); // Allocate
        memory
6      if (str == NULL) {
7          exit(1);
8      }
9      strcpy(str, "Hello, World!"); // Copy string to allocated memory
10
11     free(str); // Free the allocated memory
12
13     // Use-after-free bug: str is used after being freed
14     strcpy(str, "This will cause a use-after-free error"); // Buggy
        line
15 }
```

The generated content may differ with each run. This file, for example, contains 15 lines and UAF bug on line 14.

This demonstrates that we can effectively utilize the LLM-based synthesizer for testing our OOB detector. Rather than writing test cases by hand, we can simply let the LLM generate them for us.

In our daily lives, we regularly use LLM services such as GPT and Gemini. The number of parameters for GPT-3.5 is known to be 175B, while our selected Qwen model has 7B, a difference of approximately 25x. The scale of a 175B model makes it difficult to host on a home laptop due to GPU memory capacity.

Another significant difference is inference speed. You may have noticed that the synthesizer we built takes some time to generate tokens, even with the 256-token length

constraint. In contrast, GPT and Gemini respond almost instantly. This suggests that their inference servers are equipped with extremely powerful and expensive hardware.

### 8.4.1 Homework

**Exercise 8.1** We suggest the following topics for further study and exploration:

- **Topic #1: Generating Other Types of Bugs with the LLM**
  - Readers can modify the prompt to generate other types of bugs they are interested in.
- **Topic #2: Exploring the LLM-based Synthesizer**
  - The selected LLM is a specialized model for coding tasks. You can extend its use to synthesize code for other programming challenges, such as those in software engineering. Consider what tasks in your daily work could be synthesized and automated. For further ideas, you can also refer to state-of-the-art papers.

## 8.5 Retrospect

When I study program synthesis (Programming By Example; PBE), the area was really interested and it was fun. PBE takes pairs of input and output and grammar rules and it generates the program, especially source code. Thus, program makes program. The most popular use case of PBE is Flash Fill, which recognize the pattern of excel sheet and fills empty cells where the Flash Fill [39] inserts.

After the emergence of generative LLMs, coding tasks have become one of their major strengths. I initially thought that this would make traditional program synthesis obsolete, but I have completely reversed that opinion. The significant difference between the two is that a traditional program synthesizer yields a deterministic output for a given input (or a semantically equivalent one), whereas an LLM does not. If you ask an LLM the same prompt twice, the answer will be different, and sometimes the semantic meaning will be entirely different as well.

In this sense, I believe that guaranteeing a deterministic output is still a requirement in certain environments, such as mission-critical applications. Due to the unexplainability of an LLM's inference behavior, a traditional algorithmic approach remains attractive. I personally have concerns that the unpredictable and uncontrollable nature of LLM responses could present significant hurdles in some areas.

Therefore, it is crucial to first inspect an application's characteristics to determine if it is suitable for leveraging the power of LLMs, or if a more conservative, deterministic approach is required.



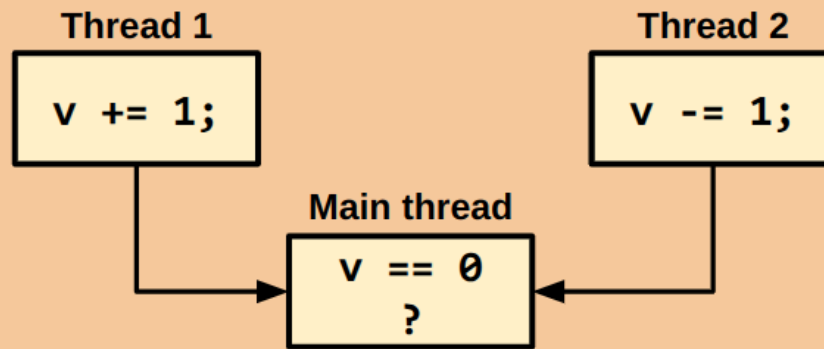
**Part VIII**

# Data Race Detector

|          |                                        |            |
|----------|----------------------------------------|------------|
| <b>9</b> | <b>Introduction to Data Race . . .</b> | <b>187</b> |
| 9.1      | Introduction . . . . .                 | 187        |
| 9.2      | Theory . . . . .                       | 191        |



# “How to detect data race?”



## 9. Introduction to Data Race

### 9.1 Introduction

Over time, hardware processor performance has steadily improved. One of the major factors driving this is the increase in the number of CPU cores. With two cores, a processor can execute two tasks simultaneously; with four cores, it can handle four tasks at once. If a processor is limited to a single core, no parallelism is possible, and all tasks must be executed sequentially.

Today, nearly all commercial products are equipped with multi-core CPUs. For example, my laptop has 32 cores. While performance is significantly enhanced by having many cores, their effective utilization presents another problem. A key challenge is a data race. **A data race is a bug where multiple threads access a shared memory location concurrently, with at least one access being a write, leading to unpredictable results.** Let's consider an example below.

```
1 // example-code/race/race-example.c
2
3 #include <stdio.h>
4 #include <pthread.h>
5 #include <unistd.h>
6
7 volatile int counter = 0;
8 int n = 20000;
9
10 void* increment_counter(void* arg) {
11     for (int i = 0; i < n; i++) {
12         counter++;
13     }
14     return NULL;
15 }
16
```

```
17 int main() {
18     pthread_t thread1, thread2;
19
20     if (pthread_create(&thread1, NULL, increment_counter, NULL) !=
21         0) {
22         perror("Error creating thread 1");
23         return 1;
24     }
25     if (pthread_create(&thread2, NULL, increment_counter, NULL) !=
26         0) {
27         perror("Error creating thread 2");
28         return 1;
29     }
30     if (pthread_join(thread1, NULL) != 0) {
31         perror("Error joining thread 1");
32         return 1;
33     }
34     if (pthread_join(thread2, NULL) != 0) {
35         perror("Error joining thread 2");
36         return 1;
37     }
38
39     printf("Expected counter value: %d\n", 2 * n);
40     printf("Actual counter value: %d\n", counter);
41
42     return 0;
43 }
```

The expected output of this example is 40000, but a data race means the program is not guaranteed to yield this result every time it runs. If you don't immediately see a different value, try running the program multiple times; you should quickly observe an incorrect output.

### Note 9.1

### Volatile keyword

Note the use of the volatile [40] keyword. We added it to this example to ensure the data race is consistently observable. Without volatile, a compiler might make optimizations—such as assuming that the value of counter won't change unexpectedly—and effectively remove the loop, thus hiding the data race. The volatile keyword prevents the compiler from making such assumptions, forcing it to read the variable from memory on every access. This keyword is typically used for variables that can be modified by factors external to the current thread, such as hardware registers or other threads.

As we have seen, data races are a common and notorious bug in multi-threaded programming. A data race often does not cause an explicit crash; instead, it corrupts the value of a shared variable. This corruption can then propagate through the program, leading to unpredictable behaviors that are notoriously difficult to debug.

Let's use a Rust example.

```

1 // example-code/race/example-race.rs
2
3 use std::thread;
4
5 static mut COUNTER: i32 = 0;
6
7 fn main() {
8     let mut handles = vec![];
9
10    for _ in 0..2 {
11        let handle = thread::spawn(|| {
12            for _ in 0..20000 {
13                COUNTER += 1;
14            }
15        });
16        handles.push(handle);
17    }
18
19    for handle in handles {
20        handle.join().unwrap();
21    }
22
23    println!("Final counter value: {}", COUNTER);
24 }
```

The semantic meaning is the same as the C example, but this program is not compilable. The following error will be produced:

```

1 error[E0133]: use of mutable static is unsafe and requires
  unsafe function or block
2 --> race-example.rs:11:17
3 |
4 11 |                 COUNTER += 1;
5    |                 ^^^^^^^ use of mutable static
6 |
7    = note: mutable statics can be mutated by multiple threads:
  aliasing violations or data races will cause undefined behavior
8
```

```

9  error[E0133]: use of mutable static is unsafe and requires unsafe
   function or block
10  --> race-example.rs:21:41
11  |
12 21 |     println!("Final counter value: {}", COUNTER);
   |                                     ^^^^^^^ use of mutable
13  |     static
   |
14  |
15  = note: mutable statics can be mutated by multiple threads:
   aliasing violations or data races will cause undefined behavior
16
17 error: aborting due to 2 previous errors
18
19 For more information about this error, try `rustc --explain E0133`.

```

A simple global variable is considered unsafe in Rust and will not compile. Similarly, some modern compilers impose constraints on language expressions to produce a safe binary at compile time. To make this code compilable, you'll need to wrap it in an unsafe block.

```

1  // example-code/race/example-race-unsafe.rs 
2
3  use std::thread;
4
5  static mut COUNTER: i32 = 0;
6
7  fn main() {
8      let mut handles = vec![];
9
10     for _ in 0..2 {
11         let handle = thread::spawn(|| unsafe {
12             for _ in 0..20000 {
13                 COUNTER += 1;
14             }
15         });
16         handles.push(handle);
17     }
18
19     for handle in handles {
20         handle.join().unwrap();
21     }
22
23     unsafe {

```

```

24     println!("Final counter value: {}", COUNTER);
25 }
26 }

```

Adding the `unsafe` keyword allows Rust to compile the program, but it shifts the responsibility for potential issues to the developer. When run, the program produces an unexpected result, not 40000.

In the previous section, we were introduced to the data race problem and saw how two compilers—C and Rust—handle it. C ignores data races at compile time, whereas Rust’s ownership system is designed to prevent them.

The goal of this chapter is to develop a data race detector based on the seminal paper by [41].

## 9.2 Theory

Several seminal papers address data race detection. We will introduce two of them, and our implementation will be based on the latter paper.

### 9.2.1 Happens-Before

Nodes in a distributed system communicate by sending and receiving messages. We will demonstrate an example of a defect that can arise in the context of event ordering during this communication.

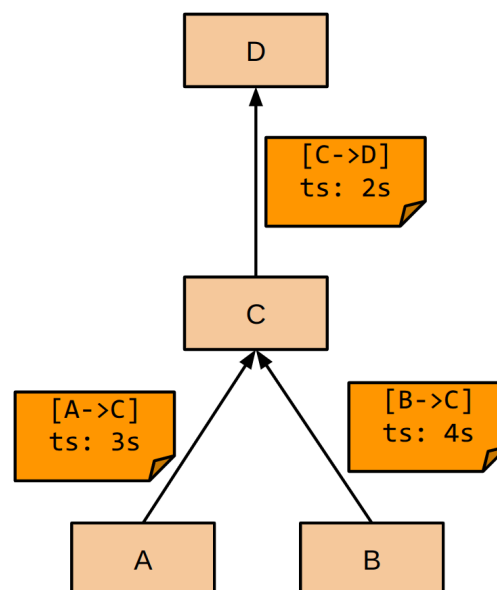


Figure 9.1: A and B send message to C, C sends to D. (ts = timestamp)

We begin with the fundamental principle that the local clocks of computers in a distributed system are never perfectly synchronized.

Consider a scenario where Node A and Node B send messages to Node C. After receiving both, Node C sends a message to Node D.

1. Node A sends its message at 3s.

2. Node B sends its message at 4s.

Let's assume Node C's local clock runs faster than Node A's and Node B's. When C receives the messages, its clock might read a time that precedes the timestamps from A and B. This mismatch reveals a critical problem: physical timestamps do not always align with the logical event order.

The logical order of events is what truly matters for causality. In our example, the sequence is as follows:

1. Node A sends a message to C, or Node B sends a message to C. The order of these two specific events doesn't matter.

2. Node C sends a message to D.

This establishes a causality between the events. Event 1 "happened-before" Event 2. This concept, known as the happened-before relation [42], considers the logical flow of events, completely ignoring the physical clock timestamps.

This is the key to synchronizing event ordering in a distributed system. Regardless of clock mismatches, each node can determine the causal relationship between events. These event relations can be chained and ordered. In our example, there are two possible causal chains:

- (1) A sends to C -> (2) B sends to C -> (3) C sends to D
- (1) B sends to C -> (2) A sends to C -> (3) C sends to D

The 'happened-before' relation can be used to identify data races. Let's look at an example below.

```
1  volatile int counter = 0;
2  int n = 20000;
3
4  void* increment_counter(void* arg) {
5      for (int i = 0; i < n; i++) {
6          counter++;
7      }
8      return NULL;
9  }
10
11 ... thread 1 and 2 execute `increment_counter()``
```

The event ordering of the increment\_counter function without a lock is as follows:

1. The local variable i is initialized to zero.
2. The global variable counter is incremented.

If we assume only a single thread exists, the counter will reach the expected value because all events are strictly ordered with no overlaps. However, in a multi-threaded environment, Event 2 can occur simultaneously across multiple threads, meaning the events are not logically ordered. This is a classic example of a data race.

Let's introduce a critical section using a mutex:

```
1  pthread_mutex_t mutex;
2  void* increment_counter(void* arg) {
3      pthread_mutex_lock(&mutex);
4      for (int i = 0; i < n; i++) {
```



```

5     counter++;
6 }
7 pthread_mutex_unlock(&mutex);
8 return NULL;
9 }

```

In a single-threaded environment, the event order is defined as:

1. The global mutex acquires a lock.
2. The local variable  $i$  is initialized to zero.
3. The global counter is incremented.
4. The global mutex releases the lock.

Just like the previous example, all events are ordered in a single-threaded context. What happens with two threads?

1. Thread 1 or Thread 2 acquires the lock on the mutex. We'll call this the "active thread." The other thread becomes the "standby thread."
2. The active thread initializes the local variable  $i$  to zero.
3. The active thread updates the global counter.
4. The active thread releases the lock.
5. The standby thread acquires the lock.
6. The standby thread initializes the local variable  $i$  to zero.
7. The standby thread updates the global counter.
8. The standby thread releases the lock.

With the critical section, all events are ordered, and no data race can occur. This demonstrates that the happened-before relation is a powerful tool for reasoning about and detecting data races.

### 9.2.2 LockSet

The paper by [41] introduces the LockSet algorithm, which we will implement. This intuitive algorithm is straightforward to understand. We've included the algorithm and a corresponding example from the paper below.

Let  $locks\_held(t)$  be the set of locks held by thread  $t$ .  
 For each  $v$ , initialize  $C(v)$  to the set of all locks.  
 On each access to  $v$  by thread  $t$ ,  
     set  $C(v) := C(v) \cap locks\_held(t)$ ;  
     if  $C(v) = \{ \}$ , then issue a warning.

Figure 9.2: Lockset Algorithm

The LockSet algorithm operates as described below. Let  $v$  be a globally shared variable, and let  $locks\_held(t)$  be the set of locks currently held by a thread  $t$ .

- **Initialization:** For each globally shared variable  $v$ , a set  $C(v)$  is initialized to contain all possible locks.
- **On Access:** Each time the variable  $v$  is accessed, the set  $C(v)$  is updated through a set intersection operation with  $locks\_held(t)$ .
- **Detection:** If the set  $C(v)$  becomes empty, a data race is detected and reported.

The following example demonstrates how the algorithm identifies a data race.

| <i>Program</i>      | <i>locks_held</i> | <i>C(v)</i>       |
|---------------------|-------------------|-------------------|
|                     | <b>{}</b>         | <b>{mu1, mu2}</b> |
| <b>lock(mu1);</b>   | <b>{mu1}</b>      |                   |
| <b>v := v+1;</b>    |                   | <b>{mu1}</b>      |
| <b>unlock(mu1);</b> | <b>{}</b>         |                   |
| <b>lock(mu2);</b>   | <b>{mu2}</b>      |                   |
| <b>v := v+1;</b>    |                   | <b>{}</b>         |
| <b>unlock(mu2);</b> | <b>{}</b>         |                   |

Figure 9.3: Example

The variable  $v$  is accessed by each thread (first three statements by thread1 and the other three statements are by thread2). The `locks_held` set is empty for the first time for each thread and  $C(v)$  is initialized with two locks (`mu1` and `mu2`). In the first thread, after `lock(mu1)` statement, `locks_held` is updated to `{mu1}`. When `v := v+1` is executed, the  $C(v)$  is refined to `{mu1}` by intersection operation as described in the algorithm. After `unlock(mu1)`, the `locks_held` becomes empty and nothing changed against  $C(v)$ .

when `lock(mu2)` is executed, `locks_held` for thread2 is updated to `{mu2}`. After `v := v+1` is executed, the  $C(v)$  is refined to empty set because the  $C(v)$  was `{mu1}`, so the set intersection yield empty set. At this point, a report is created.

Without following the algorithm step by step, we can figure out that the example contains real data race problem. The shared variable  $v$  is not protected with single mutex. Two different mutex does not make critical section for shared variable  $v$  and results in concurrent read and write can be performed against  $v$ .

If we change two statements `lock(mu2)` and `unlock(mu2)` to `lock(mu1)` and `unlock(mu1)` in the example, then no data race is reported because the intersection would yield `{mu1}` and no empty set.

Let's trace the LockSet algorithm for a shared variable  $v$ . The example code shows  $v$  being accessed by two threads.

Initially, `locks_held` is empty for both threads, and  $C(v)$  is initialized with two locks: `{mu1, mu2}`.

#### Thread 1

- `lock(mu1)`: The `locks_held` set for Thread 1 becomes `{mu1}`.
- `v := v+1`: The algorithm updates  $C(v)$  by performing an intersection:  $C(v) = C(v) \cap \text{locks\_held}(t1)$ . This refines  $C(v)$  from `{mu1, mu2}` to `{mu1}`.
- `unlock(mu1)`: The `locks_held(t1)` set becomes empty, and  $C(v)$  remains `{mu1}`.

#### Thread 2

- **lock(mu2):** The `locks_held(t2)` set for Thread 2 becomes {mu2}.
- **v := v+1:** The algorithm performs another intersection:  $C(v) = C(v) \cap \text{locks\_held}(t2)$ . The current  $C(v)$  is {mu1}, and `locks_held(t2)` is {mu2}. The intersection of these two sets is empty.  $C(v)$  becomes {}.
- **Result:** At this point, the algorithm reports a data race because  $C(v)$  is now empty.

Even without a step-by-step trace of the algorithm, we can see that this example contains a data race. The shared variable  $v$  is not protected by a single, consistent mutex. Using two different mutexes (mu1 and mu2) does not create a proper critical section for  $v$ , allowing concurrent reads and writes to occur.

If we were to change the `lock(mu2)` and `unlock(mu2)` statements to `lock(mu1)` and `unlock(mu1)`, no data race would be reported. In that case, the intersection operation would yield {mu1}, preventing  $C(v)$  from becoming empty.

To precisely detect a data race, the paper's authors introduced a more advanced algorithm: a state-transition model. The model consists of four states, as shown in the state-transition diagram:

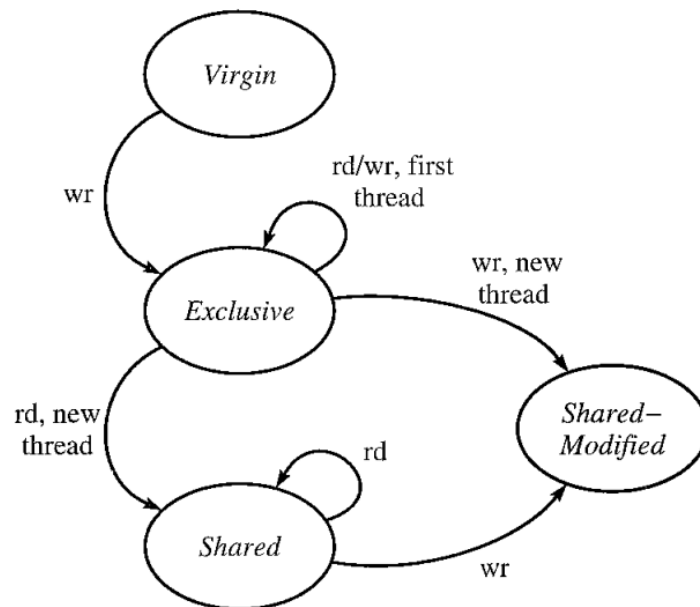


Figure 9.4: Example

- **Virgin:** This is the initial state, indicating that no operations have occurred yet.
- **Exclusive:** A state transitions to Exclusive from Virgin upon the first write operation. Subsequent read or write operations by the same thread do not change the state.
- **Shared:** A state transitions to Shared from Exclusive when another thread performs a read operation. A write operation by a different thread in this state will cause a transition to Shared-Modified.
- **Shared-Modified:** A data race is reported only in this state if the  $C(v)$  becomes empty.

Now it's time to implement this algorithm using an instrumentation module and a runtime.

### 9.2.3 Implementation

There are four tasks in the entire implementation roadmap. In our implementation, we target only static global variables. For dynamic data (heap), we will note it as homework at the end of this chapter.

- Initialization of  $C(v)$ : Implemented by static analysis
- Refinement of  $C(v)$ : Instrumented by module and runtime updates
- Update  $\text{lock\_held}(t)$ : Instrumented by module and runtime updates
- State-transaction model: Implemented by runtime

#### 9.2.3.1 Instrumentation

We introduce the entire instrumentation and static analysis loop first and explain step by step.

```

1  // instrument/src/race.rs Rust
2
3  #[derive(Debug, PartialEq, Eq, Hash, Clone)]
4  pub enum LockTyp {
5      Lock,
6      UnLock,
7  }
8  impl LockTyp {
9      pub fn i8(&self) -> i8 {
10         match self {
11             Self::Lock => 1,
12             Self::UnLock => 0,
13         }
14     }
15
16     pub fn from_i8(t: i8) -> Option<Self> {
17         match t {
18             1 => Some(Self::Lock),
19             0 => Some(Self::UnLock),
20             _ => None,
21         }
22     }
23 }
24
25 #[derive(Debug)]
26 pub enum AccessOperation {
27     Write,
28     Read,
29 }
30 impl AccessOperation {

```

```

31     pub fn i8(&self) -> i8 {
32         match self {
33             Self::Write => 1,
34             Self::Read  => 0,
35         }
36     }
37
38     pub fn from_i8(t: i8) -> Option<Self> {
39         match t {
40             1 => Some(Self::Write),
41             0 => Some(Self::Read),
42             _ => None,
43         }
44     }
45 }
46
47 #[derive(Debug, Clone)]
48 pub struct Lock<'ctx> {
49     pub typ: LockTyp,
50     pub lock: PointerValue<'ctx>,
51 }
52
53 impl<'ctx> Lock<'ctx> {
54     fn new(typ: LockTyp, lock: PointerValue<'ctx>) -> Self {
55         Self { typ, lock }
56     }
57 }
58
59 impl<'ctx> PartialEq for Lock<'ctx> {
60     fn eq(&self, other: &Self) -> bool {
61         self.typ == other.typ && self.lock.as_value_ref() ==
        other.lock.as_value_ref()
62     }
63 }
64
65 impl<'ctx> Eq for Lock<'ctx> {}
66
67 impl<'ctx> Hash for Lock<'ctx> {
68     fn hash<H: Hasher>(&self, state: &mut H) {
69         self.typ.hash(state);
70         self.lock.as_value_ref().hash(state);

```



```

111         context,
112         module,
113         builder,
114         func,
115         &instr,
116         &mut pthread_self_callsites,
117     )?;
118     build_update_lock_held(context, module,
119         builder, thread_id, lock)?;
119     continue;
120 }
121
122 if let InstructionOpcode::Load |
123     InstructionOpcode::Store = instr.get_opcode() {
124     let (operand, access_op) = match
125         instr.get_opcode() {
126             InstructionOpcode::Load => {
127                 let operand = instr.get_operand(0)
128                     .unwrap().left().unwrap();
129                 (operand, AccessOperation::Read)
130             }
131             InstructionOpcode::Store => {
132                 let operand = instr.get_operand(1)
133                     .unwrap().left().unwrap();
134                 (operand, AccessOperation::Write)
135             }
136             _ => unreachable!(),
137         };
138     let thread_id = get_thread_id(
139         context,
140         module,
141         builder,
142         func,
143         &instr,
144         &mut pthread_self_callsites,
145     )?;
146     handle_mem_access(
147         &instr,
148         &operand,
149         access_op,
150         &global_int_vals,
151         &mut lock_candidate_set,

```

```

150         &live_locks,
151         thread_id,
152         context,
153         module,
154         builder,
155     )?;
156     }
157 }
158 }
159 }
160
161 let constructor = build_race_init(context, module, builder,
162     &lock_candidate_set)?;
163 build_ctros(context, module, constructor)?;
164
165 // Verify instrumented IRs
166 module_verify(module)
167 }

```

The first task is to collect all global integer variables. (This can be extended to other types later in the homework.)

```

1 // instrument/src/race.rs Rust
2
3 fn get_global_int_vals<'ctx>(module: &Module<'ctx>) ->
4     HashSet<PointerValue<'ctx>> {
5     let globals: HashSet<PointerValue<'ctx>> = module
6         .get_globals()
7         .filter(|global| global.get_value_type().is_int_type())
8         .map(|global| global.as_pointer_value())
9         .collect();
10     globals
11 }

```

In the main loop, we collect all global variables, initialize some data structures, and iterate through the functions.

Our first task is to identify whether a given instruction is a lock or unlock operation. The `get_lock()` function checks if the instruction is a call with two operands and if the function name is either `pthread_mutex_lock` or `pthread_mutex_unlock`. The mutex variable returned from this check is then inserted into `live_locks`, which helps initialize all potential locks for the set  $C(v)$ .

Since the `lock_held(t)` function requires a thread ID, we must also instrument a call to the POSIX `pthread_self()` function. This is handled by our `get_thread_id()` function. This instrumentation is performed per each function.



```

1  // instrument/src/race.rs Rust
2
3  fn get_lock<'ctx>(instr: &InstructionValue<'ctx>) ->
   Option<Lock<'ctx>> {
4      if instr.get_opcode() != InstructionOpcode::Call ||
       instr.get_num_operands() != 2 {
5          return None;
6      }
7
8      if let Some(operand) = instr.get_operand(1) {
9          if let Some(func_name_val) = operand.left() {
10             let func_name = cstr_to_str(func_name_val.get_name());
11
12             let lock_typ = match func_name.as_str() {
13                 PTHREAD_MUTEX_LOCK => Some(LockTyp::Lock),
14                 PTHREAD_MUTEX_UNLOCK => Some(LockTyp::UnLock),
15                 _ => None,
16             };
17             if lock_typ.is_some() {
18                 let lock_operand = instr
19                     .get_operand(0)
20                     .unwrap()
21                     .left()
22                     .unwrap()
23                     .into_pointer_value();
24                 return Some(Lock::new(lock_typ.unwrap(),
25   lock_operand));
26             }
27         }
28         None
29     }
30
31 fn get_thread_id<'ctx>(
32     context: &'ctx Context,
33     module: &Module<'ctx>,
34     builder: &Builder<'ctx>,
35     func: &FunctionValue<'ctx>,
36     instr: &InstructionValue<'ctx>,
37     pthread_self_callsites: &mut HashMap<*mut LLVMValue,
38     CallSiteValue<'ctx>>,
39 ) -> Result<CallSiteValue<'ctx>> {

```

```

39     let func_id = func.as_value_ref();
40     let thread_id = match pthread_self_callsites.get(&func_id) {
41         Some(thread_id) => *thread_id,
42         None => {
43             builder.position_before(&instr);
44             let thread_id = build_pthread_self(context, module,
45                 builder)?;
46             pthread_self_callsites.insert(func.as_value_ref(),
47                 thread_id);
48         }
49     };
50     Ok(thread_id)

```

```

1  // instrument/src/inkwell_intrinsic.rs Rust
2
3  fn get_pthread_self<'ctx>(context: &'ctx Context, module:
4      &Module<'ctx>) -> FunctionValue<'ctx> {
5      match get_func(module, PTHREAD_SELF) {
6          Some(func) => func,
7          None => {
8              let pthread_self_type = context.i64_type().fn_type(&[],
9                  false);
10             module.add_function(PTHREAD_SELF, pthread_self_type,
11                 None)
12         }
13     }
14
15     pub fn build_pthread_self<'ctx>(
16         context: &'ctx Context,
17         module: &Module<'ctx>,
18         builder: &Builder<'ctx>,
19     ) -> Result<CallSiteValue<'ctx>> {
20         // instrument POSIX `pthread_self()` call
21         let pthread_self_fn = get_pthread_self(context, module);
22         let thread_id = builder.build_call(pthread_self_fn, &[], "")?;
23         Ok(thread_id)
24     }

```

We insert `lock_held(t)` right after POSIX lock and unlock calls.

```

1  // instrument/src/inkwell_intrinsic.rs Rust

```

```

2
3  fn get_update_lock_held<'ctx>(
4      context: &'ctx Context,
5      module: &Module<'ctx>,
6  ) -> FunctionValue<'ctx> {
7      match get_func(module, RACE_LOCK_HELD) {
8          Some(func) => func,
9          None => {
10             let lock_held = context.void_type().fn_type(
11                 &[
12                     context.i8_type().into(),
13                     context.i64_type().into(),
14                     context.i64_type().into(),
15                 ],
16                 false,
17             );
18             module.add_function(RACE_LOCK_HELD, lock_held, None)
19         }
20     }
21 }
22
23 pub fn build_update_lock_held<'ctx>(
24     context: &'ctx Context,
25     module: &Module<'ctx>,
26     builder: &Builder<'ctx>,
27     thread_id: CallSiteValue<'ctx>,
28     lock: Lock<'ctx>,
29 ) -> Result<()> {
30     // instrument `__race_update_lock_held()`
31     let lock_held = get_update_lock_held(context, module);
32     builder.build_call(
33         lock_held,
34         &[
35             context
36                 .i8_type()
37                 .const_int(lock.typ.i8() as u64, true)
38                 .into(),
39             thread_id.try_as_basic_value().left().unwrap().into(),
40             context
41                 .i64_type()
42                 .const_int(lock.lock.as_value_ref() as u64, true)

```

```

43         .into(),
44     ],
45     "",
46     )?;
47     Ok(())
48 }

```

The instrumented function is named `__race_update_lock_held()`. It takes three parameters:

- Lock type: An integer indicating the operation (0 for unlock, 1 for lock).
- Thread ID: Thread identifier
- ID: A mutex identifier used for reporting purposes.

```

1 %5 = call i64 @pthread_self()
2 ...
3 %6 = call i32 @pthread_mutex_lock(ptr noundef @mutex1) #6, !dbg !81
4 call void @__race_update_lock_held(i8 1, i64 %5, i64
  107913358915488),
5 ...
6 %22 = call i32 @pthread_mutex_unlock(ptr noundef @mutex1) #6
7 call void @__race_update_lock_held(i8 0, i64 %5, i64 107913358915488)
8 ...

```

LLVM-IR

In the following steps, we check if an instruction is a read or write operation. This is done by looking for load and store instructions. A load instruction indicates a read operation, and a store instruction indicates a write operation. Each of these instructions serves as a capture point to update  $C(v)$ .

This process is handled by the `handle_mem_access()` function.

```

1 // instrument/src/race.rs
2
3 fn handle_mem_access<'ctx>(<
4     instr: &InstructionValue<'ctx>,
5     operand: &BasicValueEnum<'ctx>,
6     access_op: AccessOperation,
7     global_int_vals: &HashSet<PointerValue<'ctx>>,
8     lock_candidate_set: &mut HashMap<PointerValue<'ctx>,
9     Vec<Lock<'ctx>>>,
10    live_locks: &HashSet<Lock<'ctx>>,
11    thread_id: CallSiteValue<'ctx>,
12    context: &'ctx Context,
13    module: &Module<'ctx>,
14    builder: &Builder<'ctx>,
15 ) -> Result<()> {
16     if operand.is_pointer_value() {

```

Rust

```

16         if let Some(global_var) =
17             global_int_vals.get(&operand.into_pointer_value()) {
18                 lock_candidate_set.insert(
19                     operand.into_pointer_value(),
20                     live_locks.clone().into_iter().collect::<Vec<_>>(),
21                 );
22                 if let Some(next_instr) = instr.get_next_instruction() {
23                     builder.position_before(&next_instr);
24                 }
25                 let (line, _) = get_instr_loc(instr);
26                 build_update_shared_mem(
27                     context, module, builder, access_op, thread_id,
28                     global_var, line,
29                 )?;
30             }
31     }
32     Ok(())
33 }

```

If a read or write operation does not involve a global variable, the function returns early without making any changes. If one is found, we append the `live_locks` to the initial `C(v)`. The implementation of this builder is as follows:

[illegible]

```

20     }
21 }
22 }
23
24 pub fn build_update_shared_mem<'ctx>(
25     context: &'ctx Context,
26     module: &Module<'ctx>,
27     builder: &Builder<'ctx>,
28     access_op: AccessOperation,
29     thread_id: CallSiteValue<'ctx>,
30     global_var: &PointerValue<'ctx>,
31     line: u32,
32 ) -> Result<()> {
33     let update_shared_mem = get_updated_shared_mem(context, module);
34     builder.build_call(
35         update_shared_mem,
36         &[
37             context
38                 .i8_type()
39                 .const_int(access_op.i8() as u64, true)
40                 .into(),
41             thread_id.try_as_basic_value().left().unwrap().into(),
42             context
43                 .i64_type()
44                 .const_int(global_var.as_value_ref() as u64, true)
45                 .into(),
46             context.i64_type().const_int(line as u64, true).into(),
47         ],
48         "",
49     )?;
50     Ok(())
51 }

```

There are four parameters:

1. The access operation kind, which is 0 for a load instruction and 1 for a store instruction.
2. The thread ID.
3. The global variable ID.
4. The line information where the instruction is located.

Here is the instrumented format.

```
1 %14 = load i32, ptr @counter, align 4
```

LLVM-IR

```

2  call void @_race_update_shared_mem(i8 0, i64 %5, i64
   107913358889360, i64 14)
3  %15 = add i32 %14, 1
4  call void @_race_update_shared_mem(i8 1, i64 %5, i64
   107913358889360, i64 14)
5  ...

```

Lastly, the initial  $C(v)$  is initialized via the construction function. The builder implementation is described below.

```

1  // instrument/src/inkwell_intrinsic.rs Rust
2
3  pub fn build_race_init<'ctx>(
4      context: &'ctx Context,
5      module: &Module<'ctx>,
6      builder: &Builder<'ctx>,
7      init_lock_candidate_set: &HashMap<PointerValue<'ctx>,
   Vec<Lock<'ctx>>>,
8  ) -> Result<FunctionValue<'ctx>> {
9      // __race_module_init() - initialize race module
10     let init_func_type = context.void_type().fn_type(&[], false);
11     let init_func_type_global_var = context.void_type().fn_type(
12         &[
13             // 1. global var name
14             context.ptr_type(AddressSpace::default()).into(),
15             // 2. global var decl line
16             context.i64_type().into(),
17             // 3. global var address
18             context.i64_type().into(),
19         ],
20         false,
21     );
22     let init_func_type_lock_var = context.void_type().fn_type(
23         &[
24             // 1. global var address
25             context.i64_type().into(),
26             // 2. lock var name
27             context.ptr_type(AddressSpace::default()).into(),
28             // 3. lock decl line
29             context.i64_type().into(),
30             // 4. lock address
31             context.i64_type().into(),
32         ],

```

```
33         false,
34     );
35
36     let init_fn_global_var = module.add_function(
37         RACE_INIT_CANDIDATE_LOCKSET_GLOBAL_VAR,
38         init_func_typ_global_var,
39         None,
40     );
41     let init_fn_lock_var = module.add_function(
42         RACE_INIT_CANDIDATE_LOCKSET_LOCK_VAR,
43         init_func_typ_lock_var,
44         None,
45     );
46
47     // Create a constructor function to initialize race module
48     let constructor = module.add_function(RACE_MODULE_INIT,
49         init_func_typ, None);
50     let entry = context.append_basic_block(constructor,
51         RACE_INIT_ENTRY);
52     builder.position_at_end(entry);
53
54     for (global_var, locks) in init_lock_candidate_set {
55         // initialize global variable
56         let global_var_name = get_or_build_global_string_ptr(
57             module,
58             builder,
59             &format!(
60                 "{}{}",
61                 RACE_GLOBAL_PREFIX,
62                 &cstring_to_str(global_var.get_name())
63             ),
64             );
65         let (global_var_decl_line, _) = get_instr_loc(global_var);
66         builder.build_call(
67             init_fn_global_var,
68             &[
69                 global_var_name.as_pointer_value().into(),
70                 convert_to_int_val(context,
71                     global_var_decl_line).into(),
72                 context
73             ],
74             .i64_type()
```



```

71         .const_int(global_var.as_value_ref() as u64,
72                     true)
73         .into(),
74         "",
75     )?;
76
77     // initialize lock variable
78     for lock in locks {
79         let lock_var_name = get_or_build_global_string_ptr(
80             module,
81             builder,
82             &format!(
83                 "{}{}",
84                 RACE_GLOBAL_PREFIX,
85                 &cstring(lock.lock.get_name())
86             ),
87         )?;
88         let (lock_var_decl_line, _) =
89             get_instr_loc(&lock.lock);
90         builder.build_call(
91             init_fn_lock_var,
92             &[
93                 context
94                     .i64_type()
95                     .const_int(global_var.as_value_ref() as u64, true)
96                     .into(),
97                 lock_var_name.as_pointer_value().into(),
98                 convert_to_int_val(context,
99                     lock_var_decl_line).into(),
100                 context
101                     .i64_type()
102                     .const_int(lock.lock.as_value_ref() as u64, true)
103                     .into(),
104             ],
105             "",
106         )?;
107     }
108     builder.build_return(None)?;

```

```

108     Ok(constructor)
109 }

```

Global and mutex variables are the targets for instrumentation. Line information is also attached for reporting purposes. In the working example, we have two global variables `counter` and `n`. If a mutex is found (named, `mutex1`), the mutex variables also initialized for better reporting.

```

1  define void @__race_module_init() {
2  __race_init_entry:
3      call void @__race_init_candidate_lockset_global_var(ptr
    @__race.global.counter, i64 8, i64 107913358889360)
4      call void @__race_init_candidate_lockset_lock_var(i64
    107913358889360, ptr @__race.global.mutex1, i64 5, i64
    107913358915488)
5      call void @__race_init_candidate_lockset_global_var(ptr
    @__race.global.n, i64 9, i64 107913358913936)
6      call void @__race_init_candidate_lockset_lock_var(i64
    107913358913936, ptr @__race.global.mutex1, i64 5, i64
    107913358915488)
7      ret void
8  }

```

Now we're heading to the runtime handler implementation.

### 9.2.3.2 Runtime

The runtime exposes four functions as we instrumented in the module in the previous section.

```

1  // race_runtime/src/runtime.rs
2
3  #[no_mangle]
4  pub extern "C" fn __race_init_candidate_lockset_global_var(
5      global_var_name: *const libc::c_char,
6      global_var_decl_line: u32,
7      global_var_id: i64,
8  ) {
9      init_candidate_lockset_global_var(global_var_name,
    global_var_decl_line, global_var_id);
10 }
11
12 #[no_mangle]
13 pub extern "C" fn __race_init_candidate_lockset_lock_var(
14     global_var_id: i64,
15     lock_var_name: *const libc::c_char,
16     lock_var_decl_line: u32,

```

```

17     lock_id: i64,
18 ) {
19     init_candidate_lockset_lock_var(global_var_id, lock_var_name,
20                                     lock_var_decl_line, lock_id);
21 }
22 #[no_mangle]
23 pub extern "C" fn __race_update_lock_held(is_lock: i8, thread_id:
24     i64, lock_id: i64) {
25     update_lock_held(is_lock, thread_id, lock_id);
26 }
27 #[no_mangle]
28 pub extern "C" fn __race_update_shared_mem(
29     is_write: i8,
30     thread_id: i64,
31     global_var_id: i64,
32     line: i64,
33 ) {
34     update_shared_mem(thread_id, global_var_id);
35     state_transition(is_write, thread_id, global_var_id, line);
36 }

```

The first two public functions initialize global variables (C(v)) and mutex variables. These are then managed through a singleton data structure.

```

1  // race_runtime/src/state.rs 
2
3  fn cstr_to_string(ptr: *const libc::c_char) -> String {
4      if ptr.is_null() {
5          return "".to_string();
6      }
7      unsafe { std::ffi::CStr::from_ptr(ptr).to_string_lossy().into_owned() }
8  }
9
10 pub fn init_candidate_lockset_lock_var(
11     global_var_id: i64,
12     lock_var_name: *const libc::c_char,
13     lock_var_decl_line: u32,
14     lock_id: i64,
15 ) {
16     if lock_var_name.is_null() {

```


```
17     return;
18 }
19
20 let lock_var_name = cstr_to_string(lock_var_name)
21     .strip_prefix(RACE_GLOBAL_PREFIX)
22     .unwrap()
23     .to_string();
24 lock_metadata.lock().unwrap().insert(
25     lock_id,
26     LockMetadata::new(lock_var_name, lock_var_decl_line),
27 );
28 lock_set
29     .lock()
30     .unwrap()
31     .entry(global_var_id)
32     .or_insert_with(BTreeSet::new)
33     .insert(lock_id);
34
35 init_lock_set
36     .lock()
37     .unwrap()
38     .entry(global_var_id)
39     .or_insert_with(BTreeSet::new)
40     .insert(lock_id);
41 }
42
43 pub fn init_candidate_lockset_global_var(
44     global_var_name: *const libc::c_char,
45     global_var_decl_line: u32,
46     global_var_id: i64,
47 ) {
48     if global_var_name.is_null() {
49         return;
50     }
51
52     let global_var_name = cstr_to_string(global_var_name)
53         .strip_prefix(RACE_GLOBAL_PREFIX)
54         .unwrap()
55         .to_string();
56     global_var_metadata.lock().unwrap().insert(
57         global_var_id,
```

```

58     GlobalVarMetadata::new(global_var_name,
59                             global_var_decl_line),
60 );
61 // set init state
62 // as we only consider global variables for target shared
63 // memory,
64 // directly set `Virgin` state initially
65 init_state(global_var_id);
66 }

```

```

1  // race_runtime/src/state.rs  Rust
2
3  const RACE_TEST_ENABLED: &str = "RACE_UNIT_TEST_ENABLED";
4
5  pub struct GlobalVarMetadata {
6      global_var_name: String,
7      global_var_decl_line: u32,
8  }
9  impl GlobalVarMetadata {
10     fn new(global_var_name: String, global_var_decl_line: u32) ->
11         Self {
12         Self {
13             global_var_name,
14             global_var_decl_line,
15         }
16     }
17
18     pub struct LockMetadata {
19         lock_var_name: String,
20         lock_var_decl_line: u32,
21     }
22     impl LockMetadata {
23         fn new(lock_var_name: String, lock_var_decl_line: u32) -> Self {
24             Self {
25                 lock_var_name,
26                 lock_var_decl_line,
27             }
28         }
29     }
30
31     #[derive(Eq, Hash, PartialEq)]

```

```
32 pub struct Reported {
33     thread_id: i64,
34     global_var_id: i64,
35     global_var_decl: u32,
36     global_var_used: i64,
37 }
38 impl Reported {
39     fn new(thread_id: i64, global_var_id: i64, global_var_decl: u32,
40         global_var_used: i64) -> Self {
41         Self {
42             thread_id,
43             global_var_id,
44             global_var_decl,
45             global_var_used,
46         }
47     }
48
49     #[derive(Debug, PartialEq)]
50     pub enum State {
51         Virgin,
52         Exclusive,
53         Shared,
54         SharedModified,
55     }
56
57     pub struct ThreadState {
58         thread_ids: HashSet<i64>,
59         state: State,
60     }
61     impl Default for ThreadState {
62         fn default() -> Self {
63             Self {
64                 thread_ids: HashSet::new(),
65                 state: State::Virgin,
66             }
67         }
68     }
69
70     lazy_static::lazy_static! {
71         // <global var id, metadata>
```

```

72     pub static ref global_var_metadata: Arc<Mutex<HashMap<i64,
       GlobalVarMetadata>>> = Arc::new(Mutex::new(HashMap::new()));
73     // <lock id, metadata>
74     pub static ref lock_metadata: Arc<Mutex<HashMap<i64,
       LockMetadata>>> = Arc::new(Mutex::new(HashMap::new()));
75     // <global var id, lockset>
76     pub static ref init_lock_set: Arc<Mutex<HashMap<i64,
       BTreeSet<i64>>>> = Arc::new(Mutex::new(HashMap::new()));
77     pub static ref lock_set: Arc<Mutex<HashMap<i64, BTreeSet<i64>>>>
       = Arc::new(Mutex::new(HashMap::new()));
78     // <thread id, lockset>
79     pub static ref lock_held: Arc<Mutex<HashMap<i64,
       BTreeSet<i64>>>> = Arc::new(Mutex::new(HashMap::new()));
80     // <global var id, state>
81     pub static ref state: Arc<Mutex<HashMap<i64, ThreadState>>> =
       Arc::new(Mutex::new(HashMap::new()));
82     // <set>
83     pub static ref reported: Arc<Mutex<HashSet<Reported>>> =
       Arc::new(Mutex::new(HashSet::new()));
84 }

```

The most important data structure is `ThreadState`, which holds the current state from the four possible states and a set of thread IDs.

```

1 // race_runtime/src/state.rs Rust
2
3 fn init_state(global_var_id: i64) {
4     state
5         .lock()
6         .unwrap()
7         .insert(global_var_id, ThreadState::default());
8 }

```

When a global variable is initialized, its state is also initialized to the starting point of the state-transition diagram, which is the Virgin state.

The handler of `lock_held(t)` is implemented as below.

```

1 // race_runtime/src/state.rs Rust
2
3 pub fn update_lock_held(is_lock: i8, thread_id: i64, lock_id: i64) {
4     if let Some(lock_typ) = LockTyp::from_i8(is_lock) {
5         match lock_typ {
6             LockTyp::Lock => {
7                 lock_held
8                     .lock()
9                     .unwrap()

```

```

10         .entry(thread_id)
11         .or_insert_with(BTreeSet::new)
12         .insert(lock_id);
13     }
14     LockTyp::UnLock => {
15         if let Some(set) =
16             lock_held.lock().unwrap().get_mut(&thread_id) {
17             set.remove(&lock_id);
18         }
19     }
20 }
21 }

```

The `lock_held` data structure is a map that uses thread IDs as keys and a list of lock IDs as values. When an unlock operation is performed, the corresponding lock ID is removed from the list.

The `__race_update_shared_mem` function manages both the refinement of  $C(v)$  and the state transitions. Within its body, it calls `update_shared_mem()` to refine  $C(v)$  and `state_transition()` to handle the state changes.

```

1  // race_runtime/src/state.rs 🔒 Rust
2
3  pub fn update_shared_mem(thread_id: i64, global_var_id: i64) {
4      // calc. C(v) = C(v) n lock_held(t)
5      match lock_held.lock().unwrap().get(&thread_id) {
6          Some(holds) => {
7              let mut new_lockset = BTreeSet::new();
8              if let Some(set) =
9                  lock_set.lock().unwrap().get(&global_var_id) {
10                  for l in set {
11                      if holds.contains(l) {
12                          new_lockset.insert(*l);
13                      }
14                  }
15                  lock_set.lock().unwrap().insert(global_var_id,
16                      new_lockset);
17              }
18              None => {
19                  lock_set
20                      .lock()
21                      .unwrap()
22                      .insert(global_var_id, BTreeSet::new());

```



```

22     }
23 }
24 }

```

At each load and store instruction is met, the handler is called, and  $C(v)$  is refined.

```

1  // race_runtime/src/state.rs Rust
2
3  pub fn state_transition(is_write: i8, thread_id: i64, global_var_id:
4  i64, line: i64) {
5      if let Some(access_op) = AccessOperation::from_i8(is_write) {
6          let mut s = state.lock().unwrap();
7          match access_op {
8              AccessOperation::Write => {
9                  if let Some(cur_state) = s.get_mut(&global_var_id) {
10                     match cur_state.state {
11                         State::Virgin => {
12                             cur_state.state = State::Exclusive;
13                             cur_state.thread_ids.insert(thread_id);
14                         }
15                         State::Exclusive => {
16                             if !
17                                 cur_state.thread_ids.contains(&thread_id)
18                             {
19                                 cur_state.state =
20                                     State::SharedModified;
21                                 cur_state.thread_ids.insert(thread_id);
22                             }
23                         }
24                         State::Shared => {
25                             cur_state.state = State::SharedModified;
26                             cur_state.thread_ids.insert(thread_id);
27                         }
28                         State::SharedModified => {
29                             if let Some(set) =
30                                 lock_set.lock().unwrap().get(&global_var_id)
31                             {
32                                 if set.is_empty() {
33                                     report(thread_id, global_var_id,
34   line);
35                                 }
36                             }
37                         }
38                     }
39                 }
40             }
41         }
42     }
43 }

```

```

31         }
32     }
33 }
34 AccessOperation::Read => {
35     if let Some(cur_state) = s.get_mut(&global_var_id) {
36         if cur_state.state == State::Exclusive {
37             if !
38                 cur_state.thread_ids.contains(&thread_id) {
39                 cur_state.state = State::Shared;
40                 cur_state.thread_ids.insert(thread_id);
41             }
42         }
43     }
44 }
45 }
46 }

```

For each of the four states, the conditions and corresponding state changes are implemented in a match arm. A data race is only reported in the SharedModified state. The reporting function is implemented as follows:

```

1  // race_runtime/src/state.rs 🔒 Rust
2
3  pub fn report(thread_id: i64, global_var_id: i64, line: i64) {
4      let gv_md = global_var_metadata.lock().unwrap();
5      let md = gv_md.get(&global_var_id).unwrap();
6      let report = Reported::new(thread_id, global_var_id,
7                                md.global_var_decl_line, line);
8
9      let mut r = reported.lock().unwrap();
10     if r.get(&report).is_none() {
11         println!(
12             "{}",
13             format!(
14                 "[----- Data race detected #{}
15                 -----]",
16                 r.len()
17             )
18             .red()
19             .bold()
20         );
21     }
22     if !is_test_enabled() {

```

```

20         println!("thread id      = {}", thread_id);
21     }
22     println!("variable name      = {}", md.global_var_name);
23     println!("variable decl      = {}",
24         md.global_var_decl_line);
25     println!("variable used line = {}", line);
26     println!("[related locks]");
27
28     if let Some(set) =
29         init_lock_set.lock().unwrap().get(&global_var_id) {
30         for lock_id in set {
31             let lk_md = lock_metadata.lock().unwrap();
32             let md = lk_md.get(lock_id).unwrap();
33             println!("    - lock variable name = {}",
34                 md.lock_var_name);
35             println!("    - lock variable decl = {}",
36                 md.lock_var_decl_line);
37         }
38     }
39     println!("");
40     r.insert(report);
41 }
42
43 fn is_test_enabled() -> bool {
44     matches!(env::var(RACE_TEST_ENABLED), Ok(val) if val == "1")
45 }

```

The report provides helpful information such as the variable's name, its declaration and usage locations, and any related mutex variables.

Now, let's run the initial data race example. The code is attached below.

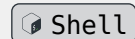
```

1  // example-code/race/race-example-fix.c
2
3  #include <stdio.h>
4  #include <pthread.h>
5  #include <unistd.h>
6
7  volatile int counter = 0;
8  int n = 20000;
9
10 void* increment_counter(void* arg) {
11     for (int i = 0; i < n; i++) {
12         counter++;

```

```
13     }
14     return NULL;
15 }
16
17 int main() {
18     pthread_t thread1, thread2;
19
20     if (pthread_create(&thread1, NULL, increment_counter, NULL) !=
21         0) {
22         perror("Error creating thread 1");
23         return 1;
24     }
25     if (pthread_create(&thread2, NULL, increment_counter, NULL) !=
26         0) {
27         perror("Error creating thread 2");
28         return 1;
29     }
30     if (pthread_join(thread1, NULL) != 0) {
31         perror("Error joining thread 1");
32         return 1;
33     }
34     if (pthread_join(thread2, NULL) != 0) {
35         perror("Error joining thread 2");
36         return 1;
37     }
38
39     printf("Expected counter value: %d\n", 2 * n);
40     printf("Actual counter value: %d\n", counter);
41
42     return 0;
43 }
```

```
1 > clang -g -O0 -S -emit-llvm race-example.c -o race-
  example.ll
2 ...
3 > ./instrument -i race-example.ll -o irace-example.ll -m race
4 ...
5 > clang irace-example.ll -L. -lrace_runtime
6 ...
7 > LD_LIBRARY_PATH=./ ./a.out
8 > LD_LIBRARY_PATH=./ ./a.out
9 [----- Data race detected #0 -----]
```



```

10 thread id      = 136849857177280
11 variable name  = counter
12 variable decl  = 5
13 variable used line = 10
14 [related locks]
15
16 [----- Data race detected #1 -----]
17 thread id      = 136849865569984
18 variable name  = counter
19 variable decl  = 5
20 variable used line = 10
21 [related locks]
22
23 Expected counter value: 40000
24 Actual counter value: 27338

```

Let's fix the example program and confirm that our detector no longer reports a data race.

```

1  #include <stdio.h>
2  #include <pthread.h>
3  #include <unistd.h>
4
5  pthread_mutex_t mutex;
6  volatile int counter = 0;
7  int n = 20000;
8
9  void* increment_counter(void* arg) {
10     pthread_mutex_lock(&mutex);
11     for (int i = 0; i < n; i++) {
12         counter++;
13     }
14     pthread_mutex_unlock(&mutex);
15     return NULL;
16 }
17
18 int main() {
19     pthread_t thread1, thread2;
20
21     if (pthread_create(&thread1, NULL, increment_counter, NULL) !=
22         0) {
23         perror("Error creating thread 1");
24         return 1;
25     }
26 }

```

```

24     }
25     if (pthread_create(&thread2, NULL, increment_counter, NULL) !=
26         0) {
27         perror("Error creating thread 2");
28         return 1;
29     }
30     if (pthread_join(thread1, NULL) != 0) {
31         perror("Error joining thread 1");
32         return 1;
33     }
34     if (pthread_join(thread2, NULL) != 0) {
35         perror("Error joining thread 2");
36         return 1;
37     }
38     printf("Expected counter value: %d\n", 2 * n);
39     printf("Actual counter value: %d\n", counter);
40
41     return 0;
42 }

```

The compile and run command is as follows:

```

1  > clang -g -O0 -S -emit-llvm race-example-fix.c -o race-
    example-fix.ll
2  ...
3  > ./instrument -i race-example-fix.ll -o irace-example-fix.ll -m race
4  ...
5  > clang irace-example-fix.ll -L. -lrace_runtime
6  ...
7  > LD_LIBRARY_PATH=./ ./a.out
8  Expected counter value: 40000
9  Actual counter value: 40000

```

A misuse of a mutex can lead to a data race. This is the last demonstration, and our detector should report it.

```

1  // example-code/race/race-example-misuse-mutex.c
2
3  #include <stdio.h>
4  #include <pthread.h>
5  #include <unistd.h>
6
7  pthread_mutex_t mutex1;

```

```
8  pthread_mutex_t mutex2;
9  volatile int counter = 0;
10 int n = 20000;
11
12 void* increment_counter1(void* arg) {
13     pthread_mutex_lock(&mutex1);
14     for (int i = 0; i < n; i++) {
15         counter++;
16     }
17     pthread_mutex_unlock(&mutex1);
18     return NULL;
19 }
20
21 void* increment_counter2(void* arg) {
22     pthread_mutex_lock(&mutex2);
23     for (int i = 0; i < n; i++) {
24         counter++;
25     }
26     pthread_mutex_unlock(&mutex2);
27     return NULL;
28 }
29
30 int main() {
31     pthread_t thread1, thread2;
32
33     if (pthread_create(&thread1, NULL, increment_counter1, NULL) !=
34         0) {
35         perror("Error creating thread 1");
36         return 1;
37     }
38     if (pthread_create(&thread2, NULL, increment_counter2, NULL) !=
39         0) {
40         perror("Error creating thread 2");
41         return 1;
42     }
43     if (pthread_join(thread1, NULL) != 0) {
44         perror("Error joining thread 1");
45         return 1;
46     }
47     if (pthread_join(thread2, NULL) != 0) {
48         perror("Error joining thread 2");
49         return 1;
50     }
```

```

48     }
49
50     printf("Expected counter value: %d\n", 2 * n);
51     printf("Actual counter value: %d\n", counter);
52
53     return 0;
54 }

```

Our detector reports it as below:

```

1  [----- Data race detected #0 -----]
2  thread id          = 131363739002560
3  variable name      = counter
4  variable decl      = 7
5  variable used line = 22
6  [related locks]
7      - lock variable name = mutex1
8      - lock variable decl = 5
9      - lock variable name = mutex2
10     - lock variable decl = 6
11
12 [----- Data race detected #1 -----]
13 thread id          = 131363747395264
14 variable name      = counter
15 variable decl      = 7
16 variable used line = 13
17 [related locks]
18     - lock variable name = mutex1
19     - lock variable decl = 5
20     - lock variable name = mutex2
21     - lock variable decl = 6
22
23 Expected counter value: 40000
24 Actual counter value: 25445

```

We have successfully implemented and demonstrated a practical data race detector. There are several areas for further improvement, which are left as homework for continued study.

### 9.2.4 Homework

**Exercise 9.1** We suggest the following topics for further study and exploration:



- **Topic #1: Extend the Race detector to capture dynamic(heap) data:** Currently, our implementation only captures static global variables. Heap data, however, can be initialized in one function and persist across other function scopes. Readers can extend the program to address this limitation.



## Part IX

# Appendix

|           |                             |            |
|-----------|-----------------------------|------------|
| <b>10</b> | <b>Conclusion . . . . .</b> | <b>229</b> |
| 10.1      | Further Study . . . . .     | 229        |



# Next Journey



## 10. Conclusion

### 10.1 Further Study

This book originated from my personal study. While I have some experience with static analysis, I had little practical experience in implementing dynamic analysis. My initial understanding was based on what I had heard about the problems and solutions in this field from coworkers.

The field of PL is both challenging and fascinating, with a vast scope. To guide your further exploration, I have outlined some key topics for continued study and research.

- **Compiler:** I agree that compilers are one of the most important components in computer science. Their applications are incredibly vast.

For example, the demand for new hardware to achieve high performance in AI systems is growing rapidly. Compilers must be able to generate efficient code for these target machines, which now include specialized hardware like ASICs. This has created a need for developing new compiler backends.

Furthermore, many concepts have been developed to generate faster code, such as operator fusion, code reordering (scheduling), and equality saturation. In the context of program analysis and engineering, a deep understanding of compiler passes—and the ability to modify them—is a valuable skill.

Some projects develop an IR from scratch for their optimization passes. To better meet the needs of new domains, multi-level IR is a popular concept for enabling more effective optimization. LLVM, for instance, allows programmers to define and generate their own custom IR and apply fine-grained, layer-specific optimizations.

- **Type System:** The type system is a fascinating and challenging area of computer science. In 2025, most programmers would agree on its importance. The success of languages like TypeScript and the growing popularity of Rust are testaments to this. A strong type system's main contribution is providing a safe programming environment. By catching type-related bugs at compile time, it eliminates a significant class of runtime errors that are often difficult to debug.

Type inference, a concept developed from type theory, allows a compiler to determine a variable's type without explicit annotation.

In new domains, such as AI compilers, there's a growing demand to design new type systems that can efficiently handle complex mathematical calculations.

- **Program Analysis:** Program bugs are ubiquitous. There are two primary approaches to finding them: static analysis and dynamic analysis. In this book, we have explored the dynamic-analysis approach, but static analysis is also widely adopted.

The key difference lies in their scope and guarantees. Static analysis, in theory, can find all bugs within a defined domain, but it often involves a trade-off between time constraints and accuracy. Achieving perfect accuracy in a reasonable amount of time is not theoretically guaranteed.

In contrast, if a bug is detected by a dynamic analysis tool, it is a real, confirmed bug. However, dynamic analysis is limited in its scope because it is impossible to explore all possible program states at runtime. This is why hybrid analysis, which combines both approaches, has been so successful in academia—it amortizes the weaknesses of each method.

Developing a strong understanding of both static and dynamic analysis is crucial. This combined knowledge provides a more robust toolkit for tackling new problems and designing effective solutions.

- **Verification:** Even with the best software hardening efforts from program analyzers, programs can still contain unknown bugs. So, how can we guarantee that a program is bug-free?

In the real world, we use unit tests to check a function's behavior across a few scenarios and inputs. For example, if you implement an `add(a, b)` function, you might write tests like `add(1, 2)` or `add(100, 200)`. However, no matter how many test cases you write, you can never achieve a theoretical guarantee of correctness.

This is where formal verification comes in. It guarantees correctness by defining theorems and lemmas about an operation, such as addition, and then formally verifying that the implementation meets those definitions. Once verified, it can generate certified code in a high-level language.

While formal verification is costly, difficult to learn, and challenging to apply broadly, it is a powerful tool for ensuring correctness in specific, mission-critical domains. It's widely used in areas like compilers, programming languages, aircraft control systems, and other applications where correctness is paramount.

- **Hardware Abstraction:** Even though hardware design and development aren't directly correlated with programming languages, the area is promising and interesting. Here are a few minor points to consider:

Abstraction is crucial in many fields, including programming language design and microarchitecture. However, the maturity of hardware abstraction lags behind the rapid growth of software. A robust, end-to-end abstraction—covering the processor pipeline, ISA, and communication with co-processors and drivers—would make the architecture more transparent and intuitive for software developers who interact with it, such as those working on compilers and SDKs.

- **AI-based ?:** Every day, a lot of news in the AI field is announced, and the results seem incredible. Sometimes I imagine developing an AI-based drone that could perform interesting tasks like delivery, observation, or running errands. These tasks are tailored for an AI-based approach.

Likewise, there are many open things to investigate and develop. Even if they're not used for software development, there are so many fun things in the world. It's great to imagine my own original ideas on how to use AI or to learn about others' ideas from papers.

This book began as a personal journey of study, so it may contain logical or implementation errors. I kindly ask for your patience and would be grateful for any corrections you may find. Your feedback is invaluable.

If you find this book to be a helpful resource, I would be deeply appreciative if you would recommend it to your friends and coworkers. Thank you so much for taking the time to read this book.

Sincerely,

Hyunsoo Shin

- 3. Review 요청
- 3. [7D] 웹개발 (gemini CLI)
- 5. [1D] avante 써보자





# Bibliography

- [1] Wikipedia, “Language-based security.”
- [2] “Clang Compiler,” [Online]. Available: <https://clang.llvm.org/>
- [3] “LLVM Pass Tutorial,” [Online]. Available: <https://llvm.org/docs/WritingAnLLVMNewPMPPass.html>
- [4] “POSIX,” [Online]. Available: <https://en.wikipedia.org/wiki/POSIX>
- [5] “Reflection Programming,” [Online]. Available: [https://en.wikipedia.org/wiki/Reflective\\_programming](https://en.wikipedia.org/wiki/Reflective_programming)
- [6] “Single Static Assignment,” [Online]. Available: [https://en.wikipedia.org/wiki/Static\\_single-assignment\\_form](https://en.wikipedia.org/wiki/Static_single-assignment_form)
- [7] “Code Coverage,” [Online]. Available: [https://en.wikipedia.org/wiki/Code\\_coverage](https://en.wikipedia.org/wiki/Code_coverage)
- [8] “Hardhat Coverage,” [Online]. Available: <https://hardhat.org/hardhat-runner/docs/guides/test-contracts>
- [9] “Inkwell,” [Online]. Available: <https://github.com/TheDan64/inkwell>
- [10] “Construction function in C,” [Online]. Available: [https://wiki.osdev.org/Calling\\_Global\\_Constructors](https://wiki.osdev.org/Calling_Global_Constructors)
- [11] “atexit,” [Online]. Available: <https://man7.org/linux/man-pages/man3/atexit.3.html>
- [12] “Mangling,” [Online]. Available: [https://en.wikipedia.org/wiki/Name\\_mangling](https://en.wikipedia.org/wiki/Name_mangling)
- [13] “Dominance Graph,” [Online]. Available: [https://en.wikipedia.org/wiki/Dominator\\_graph\\_theory](https://en.wikipedia.org/wiki/Dominator_graph_theory)
- [14] “Undefined Behavior,” [Online]. Available: [https://en.wikipedia.org/wiki/Undefined\\_behavior](https://en.wikipedia.org/wiki/Undefined_behavior)
- [15] K. Serebryany, D. Bruening, A. Potapenko, and D. Vyukov, “AddressSanitizer: A Fast Address Sanity Checker,” 2012, [Online]. Available: <https://www.usenix.org/system/files/conference/atc12/atc12-final39.pdf>
- [16] “Uninitialized Read,” [Online]. Available: [https://en.wikipedia.org/wiki/Uninitialized\\_variable](https://en.wikipedia.org/wiki/Uninitialized_variable)
- [17] “Buffer Overflow,” [Online]. Available: [https://en.wikipedia.org/wiki/Buffer\\_overflow](https://en.wikipedia.org/wiki/Buffer_overflow)
- [18] “Double Free,” [Online]. Available: [https://en.wikipedia.org/wiki/Memory\\_safety](https://en.wikipedia.org/wiki/Memory_safety)
- [19] “Null Dereference,” [Online]. Available: [https://en.wikipedia.org/wiki/Null\\_pointer](https://en.wikipedia.org/wiki/Null_pointer)
- [20] “Clang ASAN,” [Online]. Available: <https://clang.llvm.org/docs/AddressSanitizer.html>
- [21] “Fuzzing,” [Online]. Available: <https://en.wikipedia.org/wiki/Fuzzing>
- [22] “AFL,” [Online]. Available: [https://github.com/google/AFL/blob/61037103ae3722c8060ff7082994836a794f978e/docs/technical\\_details.txt](https://github.com/google/AFL/blob/61037103ae3722c8060ff7082994836a794f978e/docs/technical_details.txt)
- [23] “mmap,” [Online]. Available: <https://man7.org/linux/man-pages/man2/mmap.2.html>
- [24] “eventfd,” [Online]. Available: <https://man7.org/linux/man-pages/man2/eventfd.2.html>

- [25] “concolic testing,” [Online]. Available: [https://en.wikipedia.org/wiki/Concolic\\_testing](https://en.wikipedia.org/wiki/Concolic_testing)
- [26] “z3,” [Online]. Available: <https://github.com/Z3Prover/z3>
- [27] “Maze Solver,” [Online]. Available: <https://feliam.wordpress.com/2010/10/07/the-symbolic-maze/>
- [28] “Delta Debugging Wiki,” [Online]. Available: [https://en.wikipedia.org/wiki/Delta\\_debugging](https://en.wikipedia.org/wiki/Delta_debugging)
- [29] “Delta Debugging Paper,” [Online]. Available: <https://www.cs.purdue.edu/homes/xyzhang/fall07/Papers/delta-debugging.pdf>
- [30] “Large Language Model,” [Online]. Available: [https://en.wikipedia.org/wiki/Large\\_language\\_model](https://en.wikipedia.org/wiki/Large_language_model)
- [31] “Convolutional Neural Network,” [Online]. Available: [https://en.wikipedia.org/wiki/Convolutional\\_neural\\_network](https://en.wikipedia.org/wiki/Convolutional_neural_network)
- [32] “MNIST Dataset,” [Online]. Available: [https://en.wikipedia.org/wiki/MNIST\\_database](https://en.wikipedia.org/wiki/MNIST_database)
- [33] “Natural Language Processing,” [Online]. Available: [https://en.wikipedia.org/wiki/Natural\\_language\\_processing](https://en.wikipedia.org/wiki/Natural_language_processing)
- [34] “Long short-term memory,” [Online]. Available: [https://en.wikipedia.org/wiki/Long\\_short-term\\_memory](https://en.wikipedia.org/wiki/Long_short-term_memory)
- [35] C. S. Xia, M. Paltenghi, J. L. Tian, M. Pradel, and L. Zhang, “Fuzz4All: Universal Fuzzing with Large Language Models,” 2024, [Online]. Available: <https://arxiv.org/pdf/2308.04748>
- [36] “Fuzz4all Bug Finding,” [Online]. Available: <https://github.com/golang/go/issues/61158>
- [37] “Program Synthesis,” [Online]. Available: [https://en.wikipedia.org/wiki/Program\\_synthesis](https://en.wikipedia.org/wiki/Program_synthesis)
- [38] “Qwen Model,” [Online]. Available: <https://huggingface.co/Qwen/Qwen2.5-Coder-7B-Instruct>
- [39] “Flash Fill,” [Online]. Available: <https://support.microsoft.com/en-us/office/using-flash-fill-in-excel-3f9bcf1e-db93-4890-94a0-1578341f73f7>
- [40] “Volatile Keyword,” [Online]. Available: [https://en.wikipedia.org/wiki/Volatile\\_\(computer\\_programming\)](https://en.wikipedia.org/wiki/Volatile_(computer_programming))
- [41] S. Savage, M. Burrows, G. Nelson, P. Sobalvarro, and T. Anderson, “Eraser: a dynamic data race detector for multithreaded programs,” 1997, [Online]. Available: <https://dl.acm.org/doi/10.1145/265924.265927>
- [42] L. Lamport, “Time, clocks, and the ordering of events in a distributed system,” 1978, [Online]. Available: <https://dl.acm.org/doi/10.1145/359545.359563>