Lec 25: Fuzzing

CSED415: Computer Security Spring 2025

Seulbae Kim



Administrivia

- Lab 05 is due by the end of Friday, May 23
 - Attend office hours for help!
 - TA: Mondays and Thursdays 7-8 PM
 - Prof: Thursdays 1-2 PM

POSTECH

Administrivia

POSTECH

- Project presentations Next week
 - Each team: 15-minute presentation + 5-minute Q&A (20 minutes total)
 - Three teams will present on Tue, May 27
 - The other three teams will present on Thu, May 29
 - Presentation must include a demonstration (live or recorded)
 - All teams MUST submit their slides, code and/or binary, and report by **May 26**
 - Check PLMS assignment for details

Administrivia

POSTECH

• Final exam:

- Time: Thursday, June 5, 2:00-3:15 PM (75 minutes)
- Location: Classroom (Science Building II, Room #106)
- Format: Closed book, closed notes, no electronic devices allowed
 - Allowed: One-page (US letter- or A4-sized) double-sided handwritten cheat sheet
- Structure: 6 main questions (each may have sub-questions)
- Scope: Lectures 15-26, Labs 03-05

Program Analysis for Bug Finding



Motivation

POSTECH

- There are many bugs in the wild
 - Some bugs are security vulnerabilities that are exploitable by attackers



If we eliminate bugs, we can prevent attacks

Motivation

POSTECH

- CVE (Common Vulnerability Enumeration)
 - CVE is an indentifier assigned to publicly disclosed vulnerabilities
 - # vulnerabilities keeps increasing The attack surface is growing!



Key question

• Can we build a system that automatically finds bugs?



205722

Informal proof

POSTECH

- Define a function is_buggy
 - Input: A program
 - Output: 1 if the program has at least one bug, 0 if not

```
def is_buggy(prog):
    # test prog and return 1 or 0
```

Informal proof

POSTECH

Write a program buggy_prog

```
# buggy_prog.py
if is_buggy("buggy_prog.py"):
    return # do nothing
else:
    corrupt_memory()
    launch_root_shell()
    return
```

Self-contradictory! (Similar to the case of anti-virus)

Back to the question..

- Can we build a system that automatically finds bugs?
 - A a perfect bug-finding system cannot exist
- Therefore, we use <u>best-effort approaches</u> for <u>partial</u> bug identification
 - Bounded model checking
 - Static analysis
 - Dynamic analysis
 - etc.

Definition of "partial"

Soundness vs Completeness

- An algorithm is sound if every result it produces is in fact true
 - Every reported bug is real if algorithm is sound
 - Soundness guarantees that there is no false positive
 - A sound algorithm never misclassifies a non-bug as bug



Definition of "partial"

Soundness vs Completeness

- An algorithm is **complete** if it can derive all truths
 - Every real bug is reported if algorithm is complete
 - Completeness guarantees that there is no false negative
 - A complete algorithm never misclassifies a bug as non-bug

What algorithm identifies as bug	JS
All existing bugs (i.e., truth)	
	What algorithm identifies as bug All existing bugs (i.e., truth)

Perfect analysis

POSTECH

- Soundness vs Completeness
 - Perfect algorithm is sound and complete
 - Very challenging to achieve in practice

All existing bugs (i.e., truth) = What algorithm identifies as bugs

Metrics to evaluate a bug finding algorithm

POSTECH

• Precision, recall, and accuracy



U (all code)

- Precision: Quality of identification
 = TP / (TP + FP)
- Recall: Quantity of identification
 = TP / (FN + TP)
- Accuracy
 = (TP + TN) / U

Static vs Dynamic analysis

- Static analysis:
 - Examine program (binary or code) without running it
 - Examples:
 - Decompilation
 - Pointer analysis
 - Symbolic execution (Next topic)

- Dynamic analysis:
 - Monitor program's runtime behavior during execution
 - Examples:
 - Fuzzing (Today's topic)
 - Concolic execution





Fuzzing (or fuzz-testing)

- Definition
 - Automated testing technique that feeds invalid/unexpected/random inputs to a program under test (PUT)
 - During the process, the program is monitored for anomalous behaviors
 - Crash, hang, memory leak, etc.
 - Goal is to uncover as many bugs (and vulnerabilities) as possible

Origin of fuzzing



- Experience of Barton Miller in 1990
 - He was logged on to his workstation through a modem (dial-up line)
 - Due to a storm there were a lot of line noise (i.e., line was fuzzy)
 - The noise kept generating spurious characters on the line
 - Programs on the workstation kept crashing due to the junk characters
 - He coined the term "fuzzing" from the experience

Early days of fuzzing

• Paper: Barton Miller, et al.,

"An Empirical Study of the Reliability of Unix Utilities", Communications of the ACM, 1990



Early days of fuzzing

POSTECH

• Effectiveness

- Tested 90 Unix utility programs
 - awk, cat, cc, diff, emacs, grep, ...
- The fuzzer crashed 36 utilities!
 - Due to various bugs including unbounded pointer/array accesses, overflows, race conditions, ...
 - Randomly generated inputs were strikingly effective in triggering the bugs within poorly-written Unix programs of 1980s

- Let's put Miller's fuzzer to the test with a simple program
 - Target program reads 4 bytes from stdin
 - If the four bytes are **0xde 0xad 0xbe 0xef**, it crashes by raising segmentation fault signal

```
#include <signal.h>
#include <stdio.h>
#include <stdlib.h>
#include <unistd.h>
void bug(void) {
 printf("bug!\n");
  raise(SIGSEGV);
}
int main(void) {
 setvbuf(stdout, NULL, IONBF, 0);
  setvbuf(stdin, NULL, _IONBF, 0);
 char in[16];
  FILE *fp = fopen("/dev/stdin", "rb");
 fread(&in, 4, 1, fp);
  if (in[0] == '\xde')
    if (in[1] == '\xad')
      if (in[2] == '\xbe')
        if (in[3] == '\xef')
          bug();
 fclose(fp);
  return 0;
```

Experiment

- Let's put Miller's fuzzer to the test with a simple program
 - Fuzzer: Brute-force 4-byte random inputs until the target crashes
 - Let's check the result at the end of today's lecture

```
import os
import subprocess as sp
if __name__ == "__main__":
    trials = 0
    while True:
        __input = os.urandom(4)
        p = sp.Popen(["./target"], stdout=sp.PIPE, stdin=sp.PIPE, stderr=sp.PIPE)
        out, err = p.communicate(input=_input) # send __input to stdin and read stdout
        if b"bug!" in out:
            print(f"found in {trials} trials")
            print(f"Test input: {_input}")
            exit(0)
        print(trials)
        trials += 1
```

Interpretation of Miller's success

- Fuzzing is simple, yet effective. Why?
 - Recall the software bugs we covered in this course
 - Many attacks originate from unsanitized user inputs
 - e.g., buffer overflow, control flow hijacking, authentication bypass, DoS, SQL injection, ...
 - Fuzzing is a way to "simulate" hostile input with minimal effort

Is fuzzing still effective against modern software?

- Modern software have become very large and complex
 - Chromium browser codebase has 28 million lines of code (LoC)
 - Linux kernel comprises over 27 MLoC
 - FFmpeg has 1.4 MLoC
- Manual review of every code path is impractical
 - Imagine manually analyzing a program with the control flow graph (CFG) displayed on the right
 - Time consuming, error-prone, and hardly scalable



Is fuzzing applicable to large and complex programs?

Evolution of fuzzing

POSTECH

- Types of fuzzing
 - Blackbox, greybox, and whitebox fuzzing
 - Mutation-based vs generation-based fuzzing

Greybox Fuzzing



Overview of Black, grey, and whitebox fuzzing



- Generates random inputs
- Fuzzer has no knowledge of program's code and internal states
- The approach of Miller et al.
- Pros:
 - Extremely fast
 - Easy to use
 - Scalable
- Cons:
 - Poor effectiveness
 - Poor code coverage



- Relies on "lightweight" instrumentation of the program under test
- Fuzzer has some knowledge of the program internals during fuzzing
 - Generates semi-random inputs based on the knowledge
 - Pros: Best of both worlds
 - Scalable
 - Relatively fast
 - Decent code coverage



- Fuzzer has full knowledge of the program internals and code
- Solves path constraints to generate concrete inputs for all program branches
- Pros:
 - High code coverage
- Cons:
 - Complex
 - Slow
 - Not scalable

Breakdown of fuzzing efficiency

• A typing monkey problem

• Given infinite amount of time, can a monkey, hitting keys at random on a keyboard, type a full sentence?



The possibility is non-zero; the monkey will "almost surely" type any given sentence However, it will take astronomical amount of time

Breakdown of fuzzing efficiency

 Blackbox fuzzing **Random mutation** 0000 00**1**0 Fuzzer crash 0000 0000 Seed Test input Target system x = input() x = "LIFE"Target Seed Test input x = "LIFO" x = "5IFE" x = "LOVE" if x[0] == 'H': if x[1] == 'A': x = "HEFE" x = "DOVE" x = "LIFF"if x[2] == 'R': if x[3] == 'D': \rightarrow P(crash) = $\frac{1}{2^{32}}$ crash()

Recent breakthrough

• Greybox fuzzing with code coverage feedback



Breakdown of fuzzing efficiency

POSTECH

- A typing monkey problem (Greybox edition)
 - Keep the typed letters that are correct
 - Restart typing from the next position



Breakdown of fuzzing efficiency

POSTECH

- A typing monkey problem (Greybox edition)
 - Keep the typed letters that are correct
 - Restart typing from the next position



The possibility is dramatically increased

Coverage feedback leads to better exploration

x = "LIFE"Seed x = "5IFE" X = "LOVE"Test input X = "LIFO" $\mathbf{X} = \mathbf{HEFE}^{\mathsf{HEFE}}$ New branch. Interesting! New seed x = "HEFE" Test input x = "LEFE" x = "HAVE" New branch. Interesting! New seed x = "HAVE" → P(crash) = $\frac{1}{2^8} \times \frac{1}{4} = \frac{1}{2^{10}} > \frac{1}{2^{32}}$ Get correct byte

Select right position

POSTEEH

- Instrumentation: Modifying a program to enable analysis
 - For code coverage tracking, we want to record which branches of a program has been executed
 - We can instrument basic blocks
 - Basic block (BB): A sequence of code representing one branch of a software

- Control flow graph (CFG) of the "HARD" example
 - Consists of six basic blocks



POSTECH

Instrumentation for code coverage tracking



20572

Instrumentation for code coverage tracking



Feedback-driven greybox fuzzing is effective

POSTECH

american fuzzy lop 0.47b (readpng)		
process timing run time : 0 days, 0 hrs, 4 min, 43 a last new path : 0 days, 0 hrs, 0 min, 26 a last unig crash : none seen yet last unig hang : 0 days, 0 hrs, 1 min, 51 a	overall results sec cycles done : 0 sec total paths : 195 uniq crashes : 0 uniq hangs : 1	
cycle progress map in the state of the state of the state progress map in the state progress for the state progress is the state of th	coverage	
Total Strategy vields fuzzing strategy vields bit flips : 88/14.4k, 6/14.4k, 6/14.4k byte flips : 0/1804, 0/1786, 1/1750 arithmetics : 11/126k, 3/45.6k, 1/17.8k known ints : 11/15.8k, 4/65.8k, 6/78.2k havoc : 34/254k, 0/0 trim : 2876 8/931 (61.45% gain)	al hangs : 1 (2 unique) path geometry levels : 3 pending : 178 pending : 178 pending : 178 pending : 178 pending : 178 pending : 18 pending : 18 p	

AFL



LLVM Home | Documentation »

libFuzzer - a library for coverage-guided fuzz testing.



libFuzzer

OSS-Fuzz

Discovered millions of <u>crashes</u> in complex software systems

Test Input Generation



Mutation- vs Generation-based fuzzing

POSTECH

- Motivation: Randomly generated inputs are likely rejected by the program under test
 - e.g., When fuzzing a video player application, it is very unlikely that a fuzzer generates a properly formatted **mp4** file at random
- Two methods for better input generation
 - Mutation: Mutate a given seed to generate test inputs
 - Seed: A legitimate mp4 file
 - Generation: Generate test inputs from an input model
 - Model: Specification of **mp4** file format

Mutation

POSTECH

- Frequently used mutation operators
 - Bit-flipping: Flip a randomly selected bit
 - e.g., 0xdead (0b1101 1110 1010 $110^{0}1$) \rightarrow 0xdeaf (0b1101 1110 1010 $111^{1}1$)
 - Arithmetic operation: Select a byte and add/subtract a value
 - Randomization: Select a byte and randomize the value
 - Insertion and deletion: Add or remove bytes
 - Splicing: Crossover two test inputs
 - e.g., First half of input #1 + second half of input #2

Generation

- Generate inputs that the program under test would accept
- A model describes the correct format
 - e.g., a grammar specifying the input format
 - PNG input has header and size fields
 - The header field must have the "magic number" of PNG in order for the input to be accepted by a PNG parser



Bug Oracles



Mutation-based greybox fuzzing overview

POSTECH



A need for bug oracles

- What types of anomalous behavior do we want to find?
 - Crashes, but not all vulnerabilities lead to crashes (e.g., Lab 01)
 - Memory corruption: e.g., Use-After-Free (UAF) vulnerabilities
 - Hang: Program does not finish within a timeout period
 - Memory leaks, race conditions, specification violation, ...
- A bug oracle detects any interesting behavior occurred during the execution of a program with the test input



Bug oracles in practice

- AddressSanitizer (ASan)
 - Detects buffer overflows and use-after-free
- ThreadSanitizer (TSan)
 - Detects data races
- MemorySanitizer (MSan)
 - Detects uses of uninitialized memory

Address sanitizer

- Implemented as compiler module (available in clang and gcc)
 - Instruments all load and store instructions
 - Inserts redzones around each stack and global variable



Original program

Sanitized program

 Runtime module checks whether redzones are touched when buf is read or something is written to buf



Underflow contaminates redzone1 ASan reports buffer overflow error

Overflow contaminates redzone2 ASan reports buffer overflow error

Original program

Sanitized program

Address sanitizer in action

POSTECH

• Without ASan

```
// oob.c
#include <stdio.h>
int numbers[] = { 1, 2, 3 };
int main() { /* classic out of bounds read error. */
    printf("The 4th number in my array is: %i\n", numbers[4]);
}
```

\$ gcc oob.c -o oob

\$./oob
The 4th number in my array is: 0

The bug is missed

Address sanitizer in action

POSTECH

• With ASan

CSED415 – Spr

```
// oob.c
#include <stdio.h>
int numbers[] = { 1, 2, 3 };
int main() { /* classic out of bounds read error. */
   printf("The 4th number in my array is: %i\n", numbers[4]);
}
```

\$ gcc oob.c -fsanitize=address -o oob_asan

\$./oob_asan

Final picture

POSTECH

A coverage-based mutational greybox fuzzer



Let's check the fuzzing results (from page 23)

- How many trials were required to find the bug through blackbox fuzzing?
 - Random mutation, no coverage feedback
 - Crash: Random 4 bytes being identical to "\xde\xad\xbe\xef"
 - Theoretically requires $2^{32} \approx 4.2$ billion trials
 - Experimentally: (see terminal)



- AFL: The most widely used coverage-guided mutation-based fuzzer
 - Instrumentation for code coverage using AFL's custom complier

\$ afl-cc target.c -00 -o target_afl

• Prepare a seed input

```
$ rm -rf in out
$ mkdir in out
$ echo -ne "\xff\xff\xff\xff" > in/seed
```

• Run fuzzer

\$ afl-fuzz -i in -o out -- ./target_afl

Questions

- Is fuzzing sound? (no false positives?)
- Is fuzzing complete? (no missed bugs?)



ightarrow Fuzzer can have FP if its oracles are unsound



→ Fuzzer can miss bugs as it partially explores target program

Conclusion: Fuzzing is neither sound nor complete, but it is practical and scalable

Questions?

