# Lecture 17: Grey-box Fuzzing and Mutation Analysis

17-355/17-665/17-819: Program Analysis

Rohan Padhye

November 6, 2025

\* Course materials developed with Jonathan Aldrich and Claire Le Goues

# Puzzle: Find `x` such `p1(x)` returns `True`

```
def p1(x):
    if x * x – 10 == 15:
        return True
    return False
```

# Puzzle: Find `x` such `p2(x)` returns `True`

```python
def p2(x):
    if x > 0 and x < 1000:
        if ((x - 32) * 5/9 == 100):
            return True
    return False
```

# Puzzle: Find `x` such `p3(x)` returns `True`

```python
def p3(x):
  if x > 3 and x < 100:
    z = x - 2
    c = 0
    while z >= 2:
      if z ** (x - 1) % x == 1:
        c = c + 1
      z = z - 1
    if c == x - 3:
      return True
  return False
```
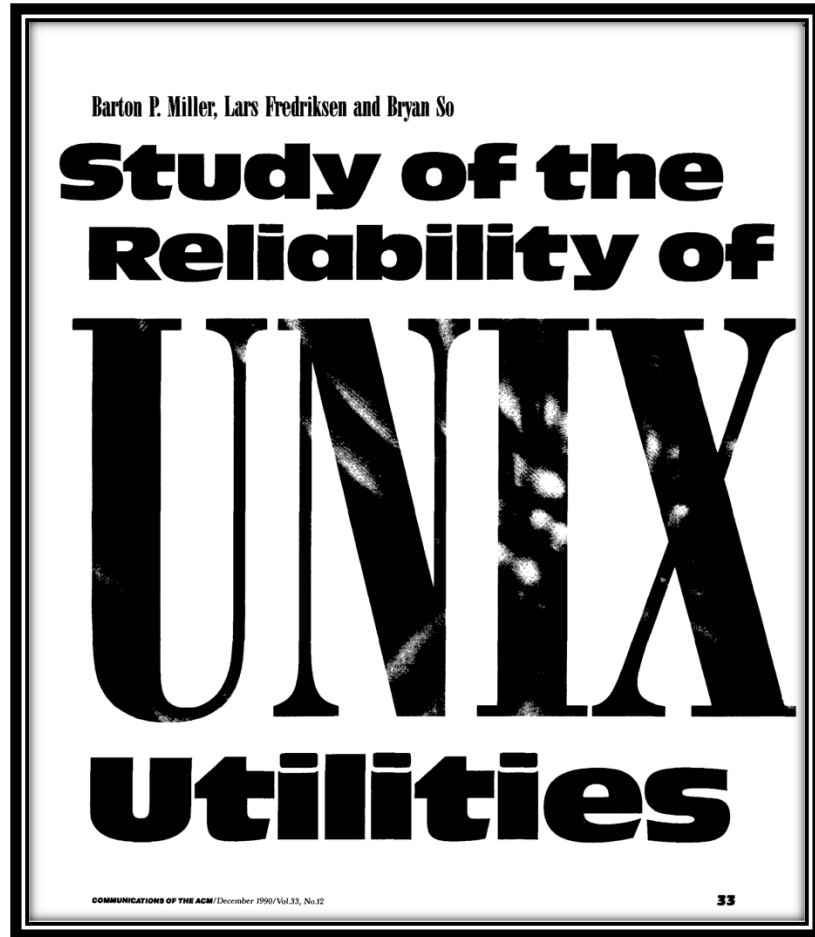
Original: https://xkcd.com/1210 CC-BY-NC 2.5

# Fuzz Testing

*Goal*:
To find program inputs that reveal a bug

*Approach:*
Generate inputs randomly until program crashes

# Fuzz Testing



**Study of the Reliability of UNIX Utilities**

Barton P. Miller, Lars Fredriksen and Bryan So

COMMUNICATIONS OF THE ACM/December 1990/Vol.33, No.12    33
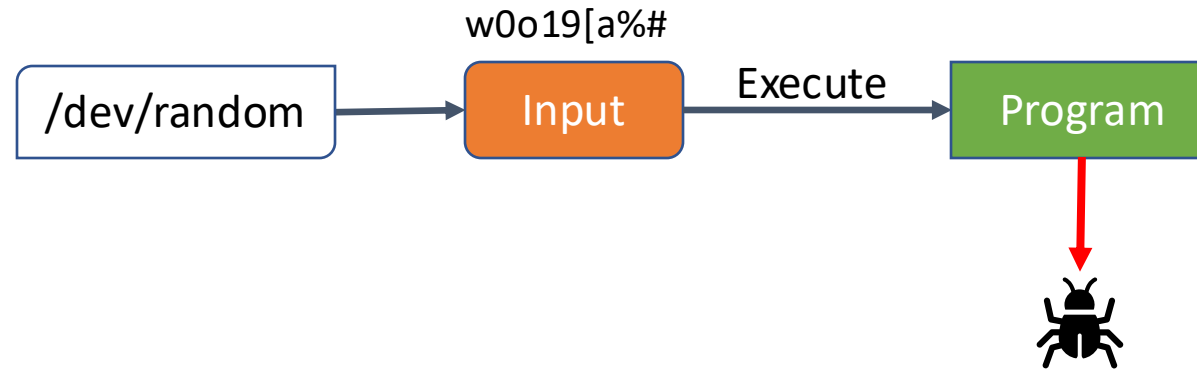
Communications of the ACM (1990)

"On a dark and stormy night one of the authors was logged on to his workstation on a dial-up line from home and the rain had affected the phone lines; there were frequent spurious characters on the line. The author had to race to see if he could type a sensible sequence of characters before the noise scrambled the command. This line noise was not surprising; but we were surprised that these spurious characters were causing programs to crash."

# Fuzz Testing 101

w0o19[a%#

/dev/random → **Input** — Execute → **Program**

🐛

1990 study found crashes in:
*adb, as, bc, cb, col, diction, emacs, eqn, ftp, indent, lex, look, m4, make, nroff, plot, prolog, ptx, refer!, spell, style, tsort, uniq, vgrind, vi*
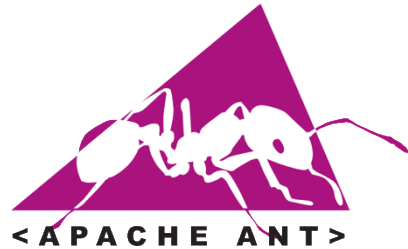
# Why do programs crash?

# Common Fuzzer-Found Bugs

<u>Causes</u>: incorrect arg validation, incorrect type casting, executing untrusted code, etc.

<u>Effects</u>: buffer-overflows, memory leak, division-by-zero, use-after-free, assertion violation, etc. ("crash")

<u>Impact</u>: security, reliability, performance, correctness

What are the **benefits**, **challenges**, & **limitations**
of this approach?

# Generate inputs randomly

$ ant –f **build.xml**

```
<project default="dist">
 <target name="init">
  <mkdir dir="${build}"/>
 </target>
 …
```

$ ant –f **/dev/random**

```
lrha3wn5p0w3uz;54 p0a23
rw3i 50a20 5a2y58a2p
y3wry3p285
q@P"uer9zparu9apur9qa3802
y5o2y 392r523a90wesu
```

Purely random data is not a very interesting input!!

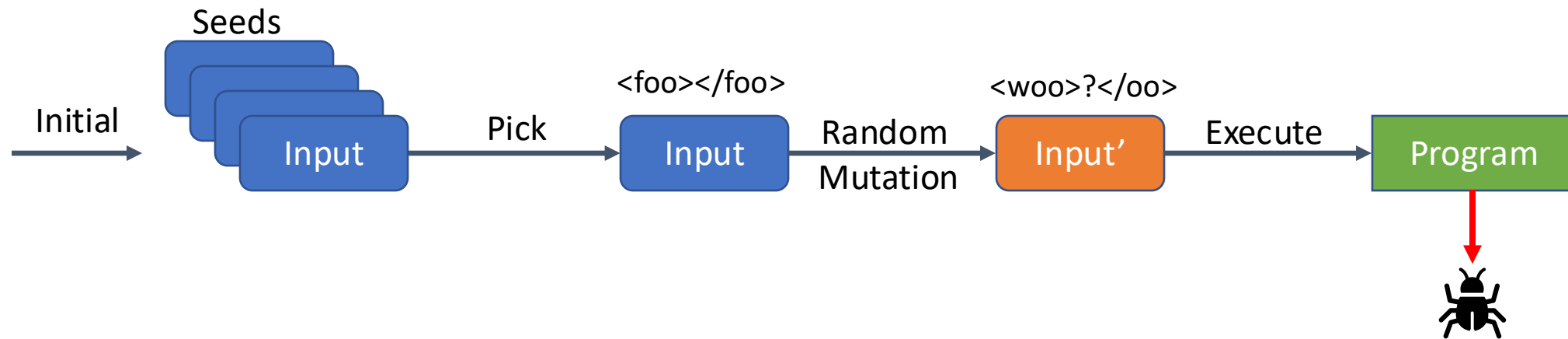# Generate inputs randomly via mutation



```
$ ant –f build.xml
```

```
<project default="dist">
 <target name="init">
  <mkdir dir="${build}"/>
 </target>
…
```

```
$ ant –f build.xml.mut
```

```
<project default="dist">
 <taWget name="init">
  <madir dir="2{build}"/@
 </tar?get>
…
```

What are some good mutations?

# Mutation Heuristics

- Binary input
  - Bit flips, byte flips
  - Change random bytes
  - Insert random byte chunks
  - Delete random byte chunks
  - Set randomly chosen byte chunks to *interesting* values e.g. INT_MAX, INT_MIN, 0, 1, -1, …
  - Other suggestions?
- Text input
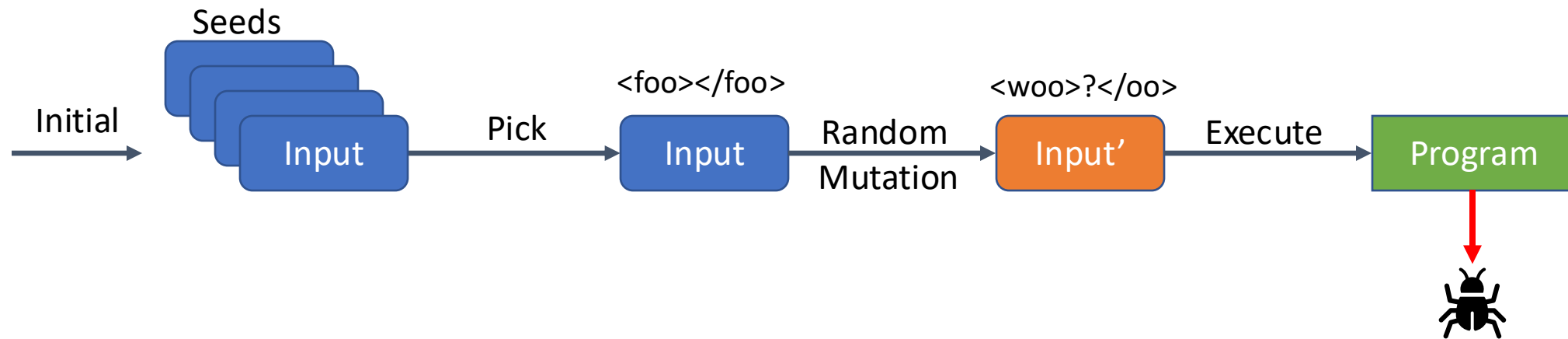  - Insert random symbols or keywords from a dictionary
  - Other suggestions?

# Mutation-Based Fuzzing (e.g. Radamsa, zzuf)

What are the **benefits**, **challenges**, & **limitations** of this approach?

How do you know if you are making progress?
Can you think of some stopping criteria?

# Code Coverage

# Exercise: How to collect coverage?

```
if (x && y) {
    s1;
    s2;
} else {
    while(b) {
        s3;
    }
}
```

# Coverage-Guided Fuzzing with AFL

# Coverage-Guided Fuzzing with AFL

November 07, 2014

## Pulling JPEGs out of thin air

This is an interesting demonstration of the capabilities of afl; I was actually pretty surprised that it worked!

```
$ mkdir in_dir
$ echo 'hello' >in_dir/hello
$ ./afl-fuzz -i in_dir -o out_dir ./jpeg-9a/djpeg
```

# Coverage-Guided Fuzzing with AFL

**The bug-o-rama trophy case**                    http://lcamtuf.coredump.cx/afl/

| | | |
|---|---|---|
| IJG jpeg [1] | libjpeg-turbo [1] [2] | libpng [1] |
| libtiff [1] [2] [3] [4] [5] | mozjpeg [1] | PHP [1] [2] [3] [4] [5] [6] [7] [8] |
| Mozilla Firefox [1] [2] [3] [4] | Internet Explorer [1] [2] [3] [4] | Apple Safari [1] |
| Adobe Flash / PCRE [1] [2] [3] [4] [5] [6] [7] | sqlite [1] [2] [3] [4] … | OpenSSL [1] [2] [3] [4] [5] [6] [7] |
| LibreOffice [1] [2] [3] [4] | poppler [1] [2] … | freetype [1] [2] |
| GnuTLS [1] | GnuPG [1] [2] [3] [4] | OpenSSH [1] [2] [3] [4] [5] |
| PuTTY [1] [2] | ntpd [1] [2] | nginx [1] [2] [3] |
| bash (post-Shellshock) [1] [2] | tcpdump [1] [2] [3] [4] [5] [6] [7] [8] [9] | JavaScriptCore [1] [2] [3] [4] |
| pdfium [1] [2] | ffmpeg [1] [2] [3] [4] [5] | libmatroska [1] |
| libarchive [1] [2] [3] [4] [5] [6] … | wireshark [1] [2] [3] | ImageMagick [1] [2] [3] [4] [5] [6] [7] [8] [9] … |
| BIND [1] [2] [3] … | QEMU [1] [2] | lcms [1] |

S3D

Carnegie
Mellon
University

# ClusterFuzz @ Chromium

# Challenging Problems

- Fuzzing heuristics
  - Mutation: Which input to mutate? How many times? Which mutations?
  - Feedback: What to instrument? How to keep overhead low?

- Oracles
  - What is a bug? Crash? Silent overflow? Infinite loop? Race condition? Undefined behavior? How do we know when we have found a bug?

- Debugging
  - Reproducibility
  - Crash triaging
  - Input minimization

- Fuzzing roadblocks
  - Magic bytes, checksums (see PNG, SSL)
  - Dependencies in binary inputs (e.g. length of chunks, indexes into tables – see PNG)
  - Inputs with complex syntax and semantics (e.g. XML, JSON, C++)
  - Stateful applications

# Oracles: Sanitizers

- Address Sanitizer (ASAN)   ***
- LeakSanitizer (comes with ASAN)
- Thread Sanitizer (TSAN)
- Undefined-behavior Sanitizer (UBSAN)

https://github.com/google/sanitizers

# AddressSanitizer https://github.com/google/sanitizers/wiki/AddressSanitizer

Compile with `clang –fsanitize=address`

Asan is a memory error detector for C/C++. It finds:

- Use after free (dangling pointer dereference)
- Heap buffer overflow
- Stack buffer overflow
- Global buffer overflow
- Use after return
- Use after scope
- Initialization order bugs
- Memory leaks

# AddressSanitizer

```
int get_element(int* a, int i) {
  return a[i];
}
```

```
int get_element(int* a, int i) {
  if (a == NULL) abort();
  return a[i];
}
```

```
int get_element(int* a, int i) {
  if (a == NULL) abort();
  region = get_allocation(a);
  if (in_stack(region)) {
   if (popped(region)) abort();
   …
  }
  if (in_heap(region)) { … }
  return a[i];
}
```

```
int get_element(int* a, int i) {
  if (a == NULL) abort();
  region = get_allocation(a);
  if (in_heap(region)) {
   low, high = get_bounds(region);
   if ((a + i) < low || (a +i) > high) {
    abort();
   }
  }
  return a[i];
}
```

# Can we go beyond coverage and crashes?

(recent-ish research results)

# Is code coverage a good measure for test-suite effectiveness?

```
bool is_even(int x) {
    int a = x / 2;
    int b = a * 2;
    return x == b;
}
```

```
void test_even() {
    assert(is_even(4) == true);
}
```

✓ 100% coverage (line, stmt, branch, path, etc.)

# *Mutation testing* measures test effectiveness on artificial bugs

```
bool is_even(int x) {
    int a = x / 2;
    int b = a * 2;
    return x == b;
}
```

```
void test_even() {
    assert(is_even(4) == true);
}
```

Kills mutant 1 & 3 but not 2

```
bool is_even(int x) {
    int a = x / 2;
    int b = a * 2;
    return x != b;
}
```
Mutant 1

```
bool is_even(int x) {
    int a = x / 2;
    int b = a * 2;
    return b == b;
}
```
Mutant 2

```
bool is_even(int x) {
    int a = x / 2;
    int b = a * 2;
    return x == a;
}
```
Mutant 3

S3D Software and Societal Systems Department

Carnegie Mellon University

# *Mutation testing* measures test effectiveness on artificial bugs

```
bool is_even(int x) {
    int a = x / 2;
    int b = a * 2;
    return x == b;
}
```

```
void test_even() {
    assert(is_even(4) == true);
    assert(is_even(1) == false);
}
```

Kills mutants 1--3

```
bool is_even(int x) {
    int a = x / 2;
    int b = a * 2;
    return x != b;
}
```

Mutant 1

```
bool is_even(int x) {
    int a = x / 2;
    int b = a * 2;
    return b == b;
}
```

Mutant 2

```
bool is_even(int x) {
    int a = x / 2;
    int b = a * 2;
    return x == a;
}
```

Mutant 3

# *Mutation testing* measures test effectiveness on artificial bugs

```
bool is_even(int x) {
    int a = x / 2;
    int b = a * 2;
    return x == b;
}
```

```
void test_even() {
    assert(is_even(4) == true);
    assert(is_even(1) == false);
}
```

Does not kill mutant 4!

```
bool is_even(int x) {
    int a = x / 2;
    int b = a + 2;
    return x == b;
}
```

Mutant 4

S3D Software and Societal Systems Department

Carnegie Mellon University

# *Mutation testing* measures test effectiveness on artificial bugs

```
bool is_even(int x) {
    int a = x / 2;
    int b = a * 2;
    return x == b;
}
```

```
void test_even() {
    assert(is_even(4) == true);
    assert(is_even(1) == false);
    assert(is_even(2) == true);
}
```

Kills mutants 1--4

```
bool is_even(int x) {
    int a = x / 2;
    int b = a + 2;
    return x == b;
}
```

Mutant 4

# *Mutation testing* measures test effectiveness on artificial bugs

```
bool is_even(int x) {
    int a = x / 2;
    int b = a * 2;
    return x == b;
}
```

```
void test_even() {
    assert(is_even(4) == true);
    assert(is_even(1) == false);
    assert(is_even(2) == true);
}
```

Does not kill mutant 5

```
bool is_even(int x) {
    int a = x / 2;
    int b = a + 2;
    return x == b;
}
```

Mutant 4

```
bool is_even(int x) {
    int a = x / 2;
    int b = a + a;
    return x == b;
}
```

Mutant 5

# *Mutation testing* measures test effectiveness on artificial bugs

```
bool is_even(int x) {
    int a = x / 2;
    int b = a * 2;
    return x == b;
}
```

**Impossible to kill mutant 5!!!**

(Mutant 5 is equivalent to original program)
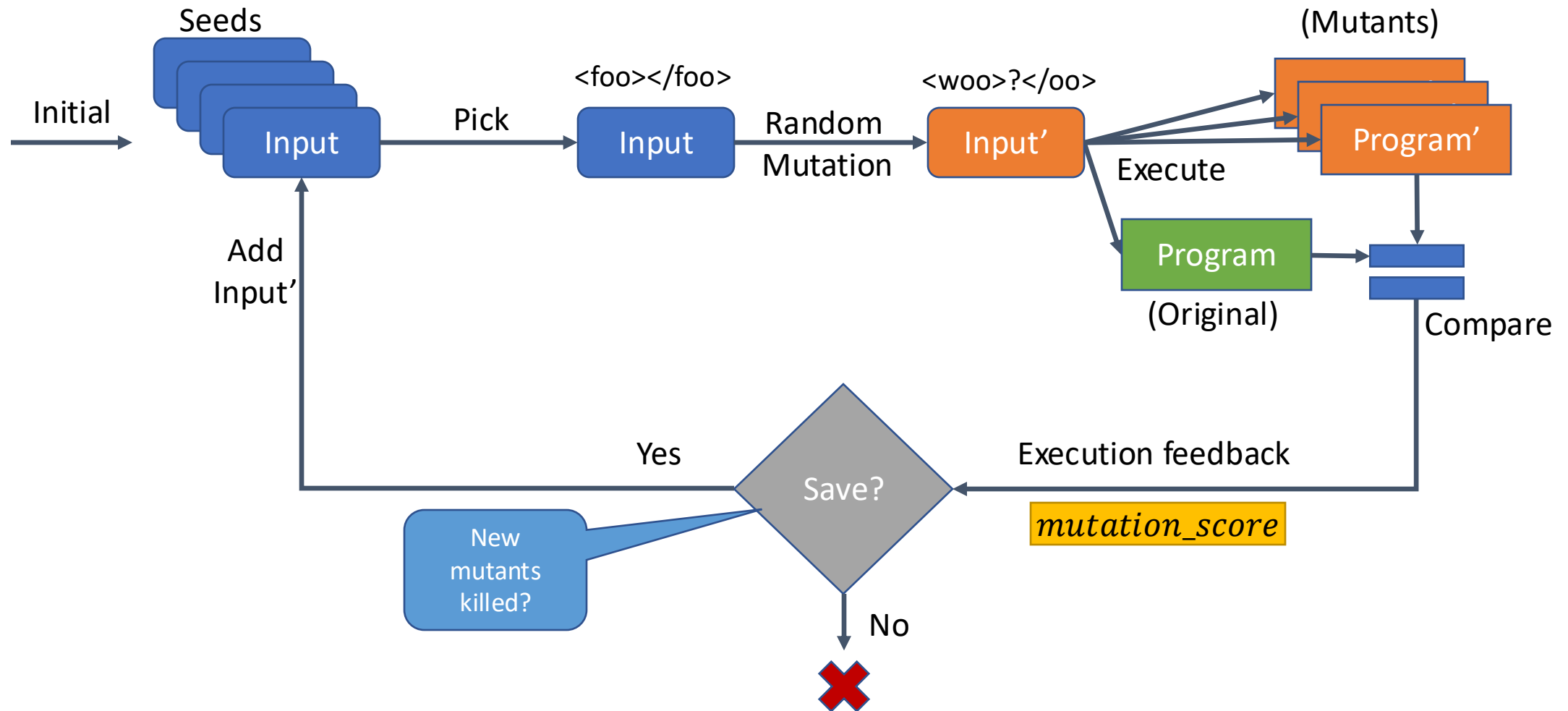
```
bool is_even(int x) {
    int a = x / 2;
    int b = a + a;
    return x == b;
}
```
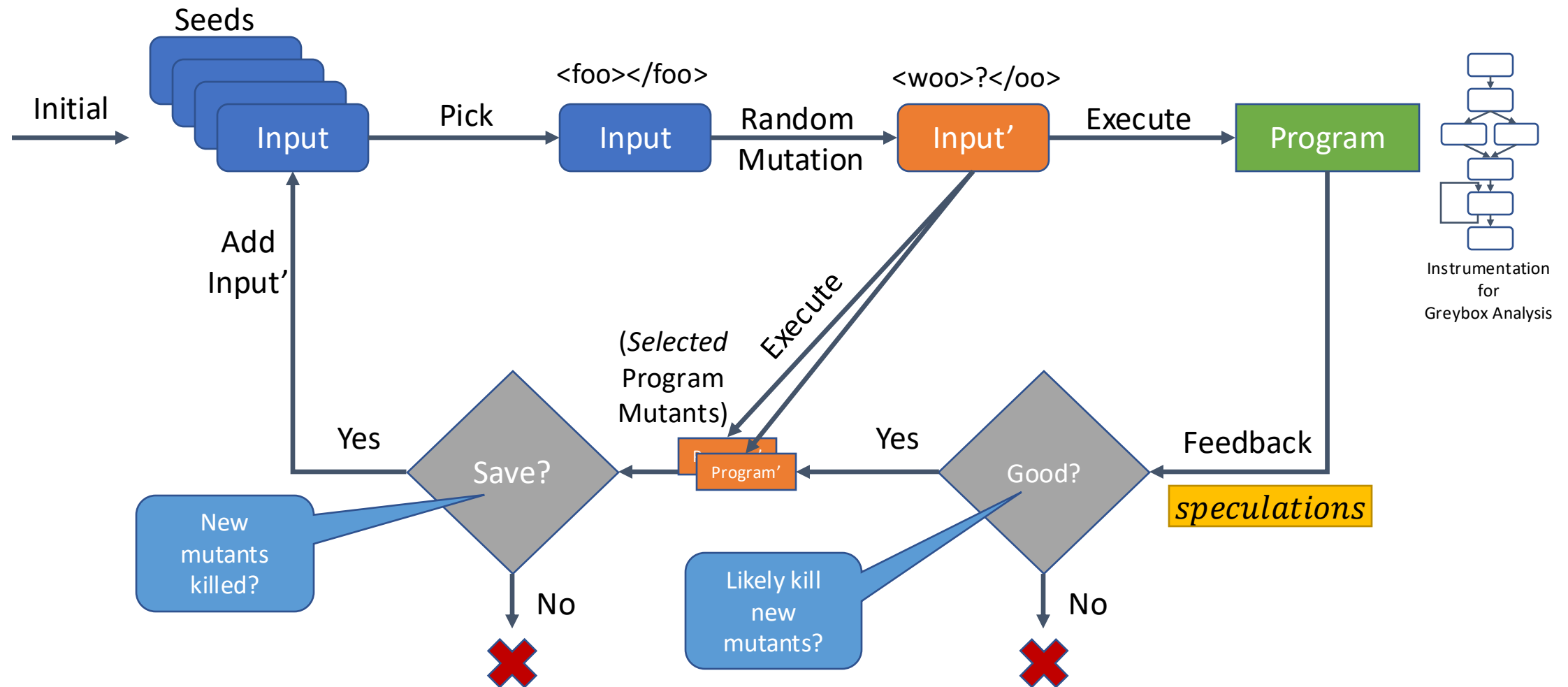
Mutant 5

# New Idea for Fuzz Testing: Mu2 – Mutation-based Mutation-guided Grey-box Fuzzing

Mutates *both* the inputs and the program

# Mutation-analysis guided fuzzing

# Speculative mutation analysis



Seeds

Initial

Input

Pick

`<foo></foo>`

Input

Random Mutation

`<woo>?</oo>`

Input'

Execute

Program

Instrumentation for Greybox Analysis

Add Input'

Execute

(*Selected* Program Mutants)

Program'

Save?

Yes

Good?

Yes

Feedback

*speculations*

New mutants killed?

No ✖

Likely kill new mutants?

No ✖

S3D — Software and Societal Systems Department

Carnegie Mellon University

# Speculative mutation analysis (PIE model)

- Let P be a program such that it's output is P(X) = Y
- Let mutant P' be "change `z: = a + b` to `z := a – b` at Line 42"
- For a given fuzzer-generated input X,
  - If X does not cover line 42, the mutant cannot be killed [Execution]
  - If X executes line 42 but in all cases `a + b == a – b` (e.g., say `b=0`) then the mutant cannot be killed [Infection]
  - If X executes line 42 but in all cases the way `z` is used does not change (e.g., the program only checks if `z > k` and `b > 0`) then the mutant cannot be killed [Propagation]
- If either of P-I-E analyses tell us mutant is unkillable, skip it!
- Great savings in fuzzing efficiency! (see ISSTA'23 paper for details)