

# Machine Programming

Lecture 15 – LLM for Software Testing

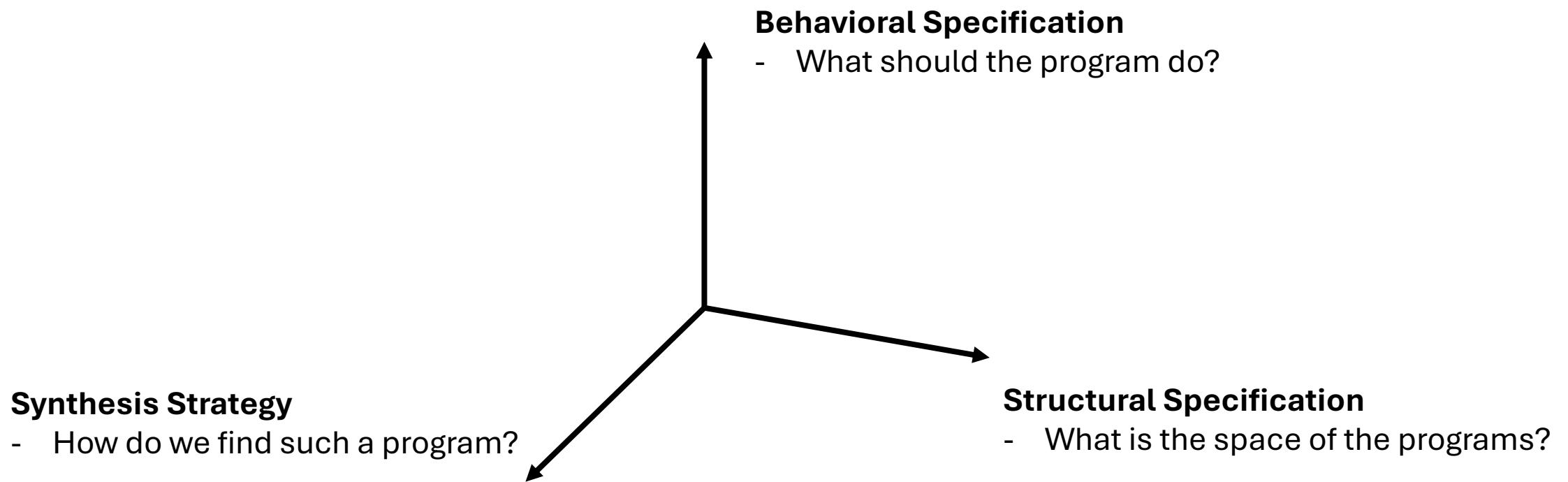
# Logistics – Week 9

- Assignment 3: Coding LLM Agents
  - <https://github.com/machine-programming/assignment-3>
  - Fully functional web-app agent. Due: Oct 23 (Thu)
- Forming groups for your final projects!
  - Form a group of 2-3 before This Thursday (~~Oct 19~~ Oct 23)

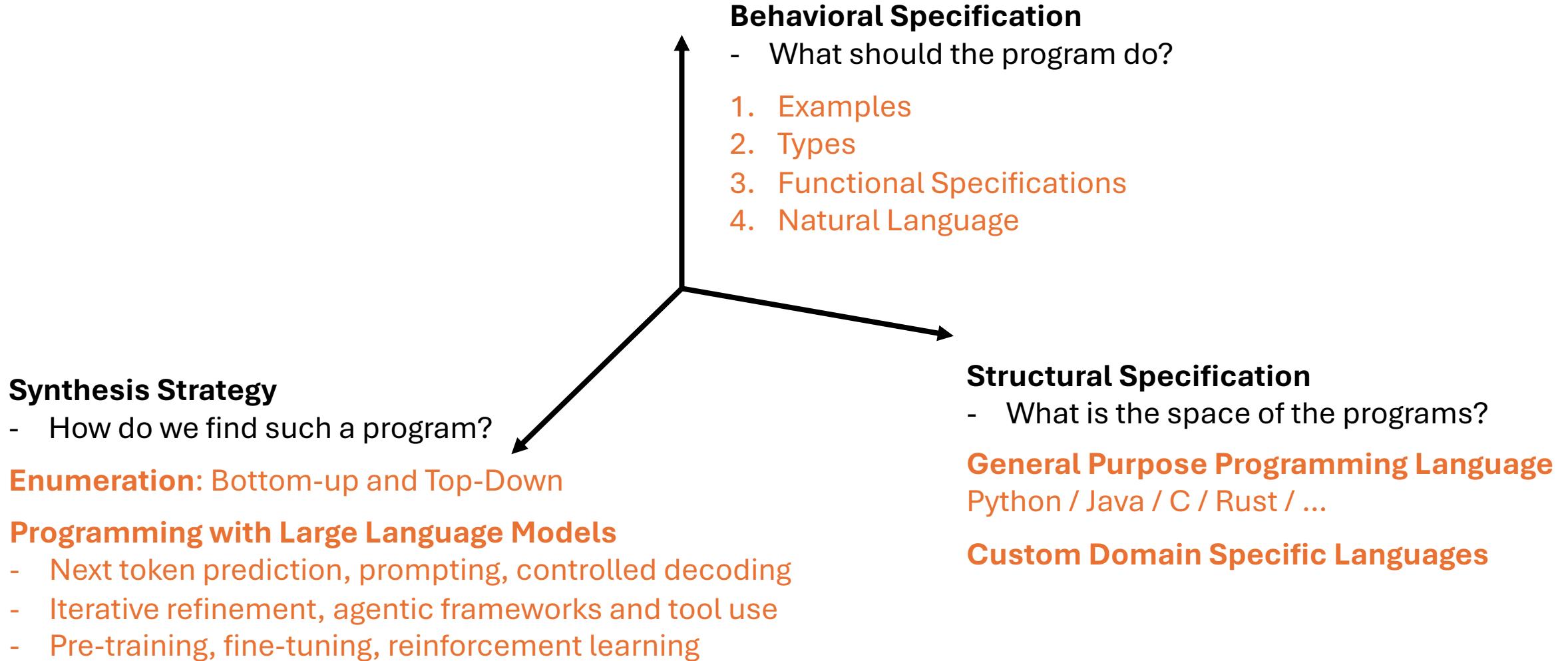
# Review

- What we have learned so far
  - Programming languages and synthesis: syntax, semantics, enumeration
  - Programming with LLMs: language modeling, prompting, agents, training

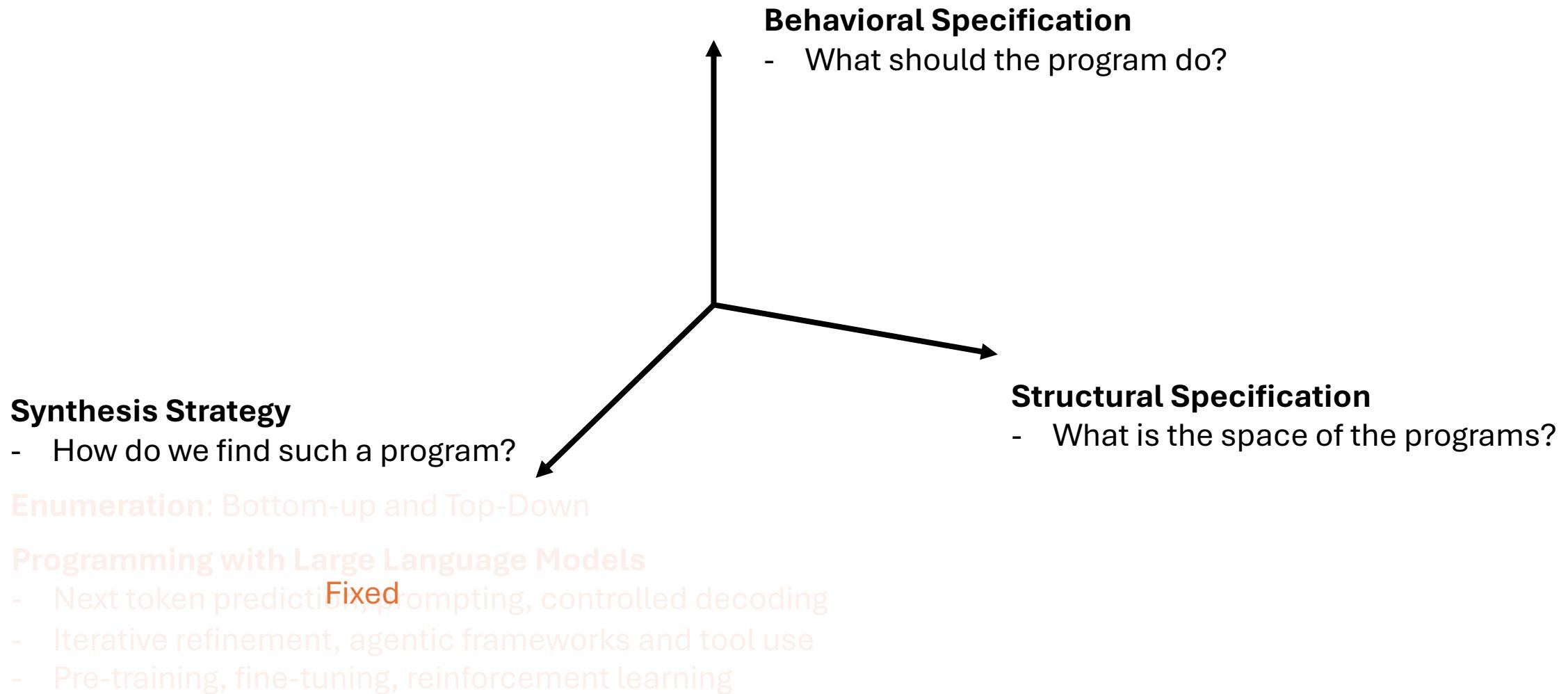
# Dimensions in Program Synthesis



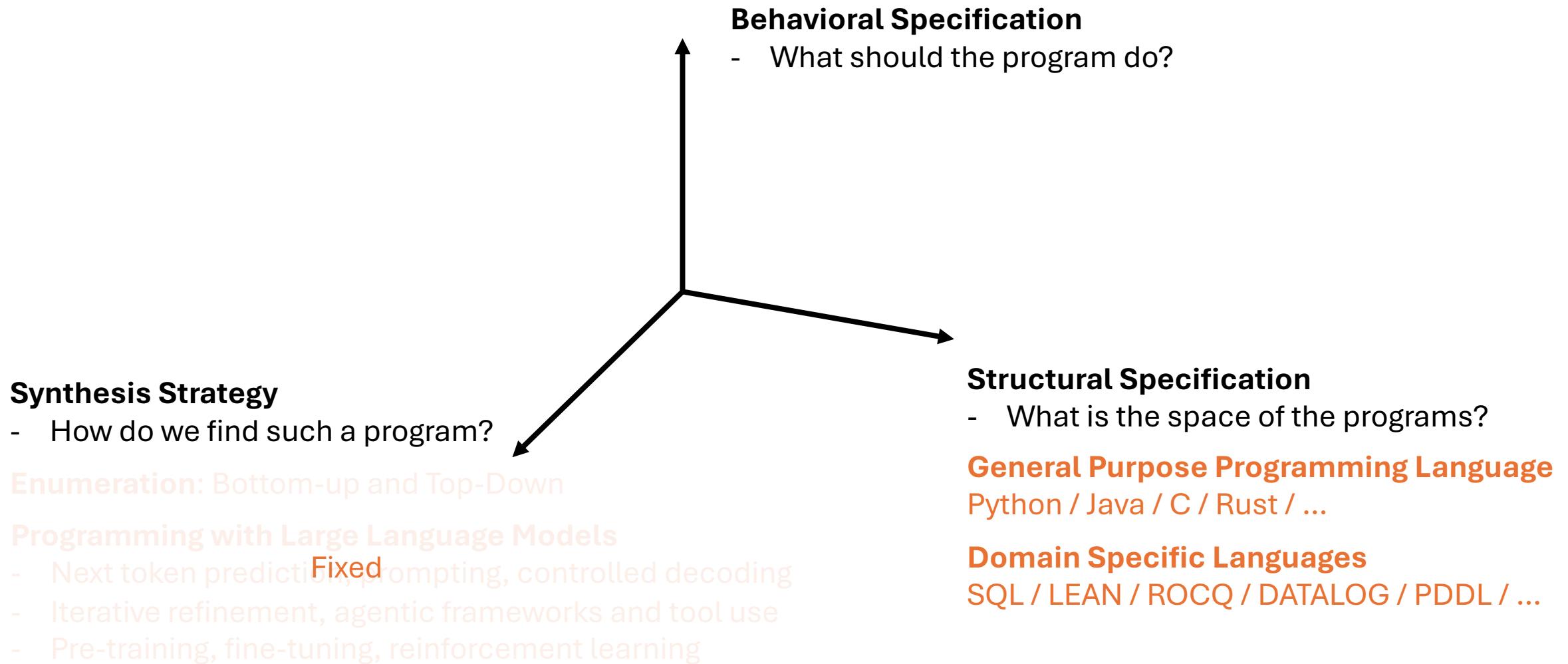
# The Course So Far



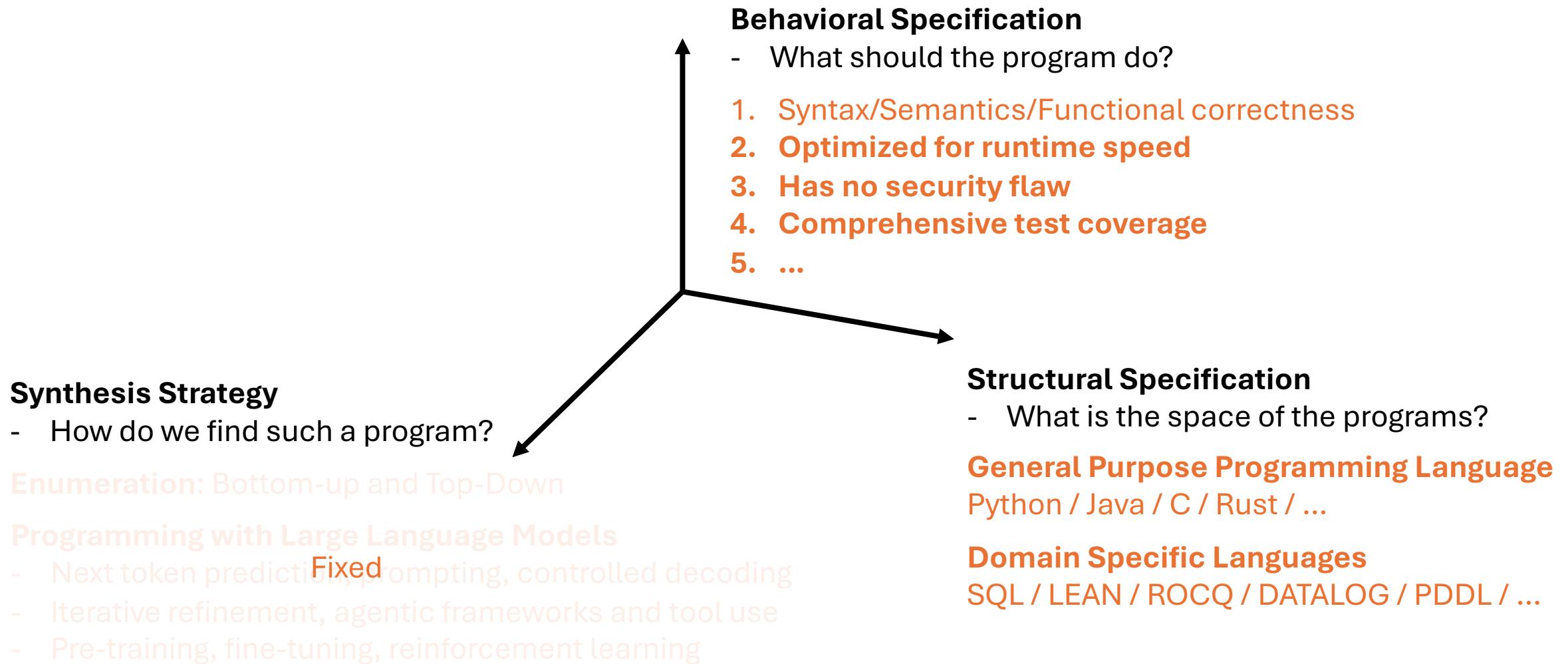
# Module 3



# Module 3



# Module 3



# Module 3: Overview

- What we have learned so far
  - Programming languages and synthesis: syntax, semantics, enumeration
  - Programming with LLMs: language modeling, prompting, agents, training
- Module 3: Applications of Machine Programming
  - Goal: how synthesis can help with diverse applications
  - Application: Software analysis, testing, optimization, security, verification
  - Application: Logic and Formal Methods, Math, Physics, Science
  - Application: Embodied systems, robotics, planning, reward modeling
  - Application: UI, interactive environment, games

# Software Analysis

Testing and Dynamic Analysis

```
BH_I16 := Integer_16(BH_F64);
```

```
BH_I16 := Integer_16(BH_F64);
```

Ada: convert 64-bit floating point number to 16-bit signed integer

# The Ariane Rocket Disaster



**Date:** June 4, 1996

**Mission:** Ariane 5 Flight 501

- European Space Agency

# The Ariane Rocket Disaster: Post Mortem

- Outcome:
  - Rocket veered off course after 37 seconds, self-destructed.

# The Ariane Rocket Disaster: Post Mortem

- Outcome:
  - Rocket veered off course after 37 seconds, self-destructed.
- What happened:
  - The rocket reused Ariane 4's Inertial Reference System (SRI) software.
  - During ascent, the SRI computed a variable BH (Horizontal Bias), representing horizontal velocity.
  - The program executed:

```
BH_I16 := Integer_16(BH_F64);
```

# The Ariane Rocket Disaster: Post Mortem

- Outcome:
  - Rocket veered off course after 37 seconds, self-destructed.
- What happened:
  - The rocket reused Ariane 4's Inertial Reference System (SRI) software.
  - During ascent, the SRI computed a variable BH (Horizontal Bias), representing horizontal velocity.
  - The program executed:

```
BH_I16 := Integer_16(BH_F64);
```

- On Ariane 4, `BH_F64` never exceeded 32767.
- On Ariane 5 (more powerful),  $\text{BH\_F64} \approx 65535$ , causing a numeric overflow.
- In Ada, this raised a `Constraint_Error` exception (like an “Operand Error”).

# The Ariane Rocket Disaster: Post Mortem

- Outcome:
  - Rocket veered off course after 37 seconds, self-destructed.
  - Cost: \$350M payload + \$150M rocket
- Have we programmers learned from the mistake?



# Common Weakness Enumeration

A community-developed list of SW & HW weaknesses that can become vulnerabilities



New to CWE?  
[Start here!](#)

Home > CWE Top 25 > 2024

ID Lookup:  Go

Home | About ▾ | Learn ▾ | Access Content ▾ | Community ▾ | Search ▾

## 2024 CWE Top 25 Most Dangerous Software Weaknesses

[Top 25 Home](#)

Share via:

[View in table format](#)

[Key Insights](#)

[Methodology](#)

1

Improper Neutralization of Input During Web Page Generation ('Cross-site Scripting')

[CWE-79](#) | CVEs in KEV: 3 | Rank Last Year: 2 (up 1)

2

Out-of-bounds Write

[CWE-787](#) | CVEs in KEV: 18 | Rank Last Year: 1 (down 1)

3

Improper Neutralization of Special Elements used in an SQL Command ('SQL Injection')

[CWE-89](#) | CVEs in KEV: 4 | Rank Last Year: 3

...

22

Use of Hard-coded Credentials

[CWE-798](#) | CVEs in KEV: 2 | Rank Last Year: 18 (down 4)

23

Integer Overflow or Wraparound

[CWE-190](#) | CVEs in KEV: 3 | Rank Last Year: 14 (down 9)

24

Uncontrolled Resource Consumption

[CWE-400](#) | CVEs in KEV: 0 | Rank Last Year: 37 (up 13)



## 2024 CWE Top 25 Most Dangerous Software Weaknesses

[Top 25 Home](#)

Share via:

[View in table format](#)

[Key Insights](#)

[Methodology](#)

1

Improper Neutralization of Input During Web Page Generation ('Cross-site Scripting')

[CWE-79](#) | CVEs in KEV: 3 | Rank Last Year: 2 (up 1)

2

Out-of-bounds Write

[CWE-787](#) | CVEs in KEV: 18 | Rank Last Year: 1 (down 1)

3

Improper Neutralization of Special Elements used in an SQL Command ('SQL Injection')

[CWE-89](#) | CVEs in KEV: 4 | Rank Last Year: 3

...

22

Use of Hard-coded Credentials

[CWE-798](#) | CVEs in KEV: 2 | Rank Last Year: 18 (down 4)

23

Integer Overflow or Wraparound

[CWE-190](#) | CVEs in KEV: 3 | Rank Last Year: 14 (down 9)

24

Uncontrolled Resource Consumption

[CWE-400](#) | CVEs in KEV: 0 | Rank Last Year: 37 (up 13)

# Goals of Software Analysis

- Question:
  - If Ariane 5's code **compiled** cleanly, why did it still **explode**?
- Idea:
  - Software analysis ≠ compile without error/warning
  - Search for hidden **failure modes** before deployment

# Goals of Software Analysis

Goal	Example	Early Warning Signal
Detect Bugs	Off-by-one, numeric overflow, null dereference, out-of-bounds	Static range warnings, failed test cases
Detect Vulnerabilities	SQL-injection, buffer overflow, path traversal, incorrect permission	CWE-based pattern matching
Check Correctness	Functional specifications, Algorithmic invariants	Violated assertions
Verification	Know whether a functional specification “always holds”	Verifier error, symbolic execution traces
Measure Reliability	Coverage (line, block, path), test pass rate	Uncovered lines, failed test cases

# Forms of Software Analysis

Static Analysis

Dynamic Analysis

# Forms of Software Analysis

## Static Analysis

Examine code without concrete execution

## Dynamic Analysis

Observes program behavior while executing

# Forms of Software Analysis

## Static Analysis

Examine code without concrete execution

- Taint analysis
- Reachability analysis
- Abstract interpretation
- Symbolic execution
- ...

## Dynamic Analysis

Observes program behavior while executing

- Unit testing
- Fuzzing
- Property-based testing
- Penetration testing
- ...

# Forms of Software Analysis

## Dynamic Analysis

(This Lesson)

Observes program behavior while executing

- Unit testing
- Fuzzing
- Property-based testing
- Penetration testing

# Dynamic Analysis

- Some bugs only manifest when code runs
  - Runtime errors
  - Unhandled exceptions
  - Errors dependent on external environments
- Want to concretely execute code
  - Observe program behavior
  - Analyze potential mistakes

# Dynamic Testing

## Unit Testing

Human-written tests  
in the form of concrete  
input-output

```
assert f(2) == 4
```

# Dynamic Testing

Unit Testing

Fuzzing

Human-written tests  
in the form of concrete  
input-output

Randomly  
generating  
test inputs

```
assert f(2) == 4          for i in 0..10^4:  
                           x = gen_input()  
                           f(x)
```

# Dynamic Testing

## Unit Testing

Human-written tests  
in the form of concrete  
input-output

```
assert f(2) == 4
```

## Fuzzing

Randomly  
generating  
test inputs

```
for i in 0..10^4:  
    x = gen_input()  
    f(x)
```

## Property-based testing

Randomly generating  
test inputs, but with  
properties checking outputs

```
for i in 0..10^4:  
    x = gen_input()  
    y = f(x)  
    assert property(x,y)
```

# Dynamic Testing

## Unit Testing

Human-written tests  
in the form of concrete  
input-output

```
assert f(2) == 4
```

## Fuzzing

Randomly  
generating  
test inputs

```
for i in 0..10^4:  
    x = gen_input()  
    f(x)
```

## Property-based testing

Randomly generating  
test inputs, but with  
properties checking outputs

```
for i in 0..10^4:  
    x = gen_input()  
    y = f(x)  
    assert property(x,y)
```

## Penetration testing

Crafting adversarial inputs  
to trigger vulnerabilities rather  
than getting outputs

```
try:  
    f(very_bad_input)  
    assert false  
except:  
    assert no_bad_behavior
```

## Unit Testing

Human-written tests  
in the form of concrete  
input-output

```
assert f(2) == 4
```

```
for i in 0..10^4:  
    x = gen_input()  
    f(x)
```

Human Effort: figure out **input-outputs**, write the **test cases**

## Fuzzing

Randomly  
generating  
test inputs

## Property-based testing

Randomly generating  
test inputs, but with  
properties checking outputs

```
for i in 0..10^4:  
    x = gen_input()  
    y = f(x)  
    assert property(x,y)
```

## Penetration testing

Crafting adversarial inputs  
to trigger vulnerabilities rather  
than getting outputs

```
try:  
    f(very_bad_input)  
    assert false  
except:  
    assert no_bad_behavior
```

Human Effort: write the **input generator** and **fuzzing harness**

Unit Testing

Fuzzing

Property-based testing

Penetration testing

Human-written tests  
in the form of concrete  
input-output

```
assert f(2) == 4
```

```
for i in 0..10^4:  
    x = gen_input()  
    f(x)
```

Randomly  
generating  
test inputs

Randomly generating  
test inputs, but with  
properties checking outputs

```
for i in 0..10^4:  
    x = gen_input()  
    y = f(x)  
    assert property(x,y)
```

Crafting adversarial inputs  
to trigger vulnerabilities rather  
than getting outputs

```
try:  
    f(very_bad_input)  
    assert false  
except:  
    assert no_bad_behavior
```

Human Effort: figure out **input-outputs**, write the **test cases**

Human Effort: write the **input generator** and **fuzzing harness**

Unit Testing

Human-written tests  
in the form of concrete  
input-output

```
assert f(2) == 4
```

```
for i in 0..10^4:  
    x = gen_input()  
    f(x)
```

Human Effort: figure out **input-outputs**, write the **test cases**

Fuzzing

Randomly  
generating  
test inputs

Property-based testing

Randomly generating  
test inputs, but with  
properties checking outputs

```
for i in 0..10^4:  
    x = gen_input()  
    y = f(x)  
    assert property(x,y)
```

Human Effort: write the **input generator**, **harness** and the **property**

Penetration testing

Crafting adversarial inputs  
to trigger vulnerabilities rather  
than getting outputs

```
try:  
    f(very_bad_input)  
    assert false  
except:  
    assert no_bad_behavior
```

Human Effort: write the **input generator** and **fuzzing harness**

Unit Testing

Human-written tests  
in the form of concrete  
input-output

```
assert f(2) == 4
```

```
for i in 0..10^4:  
    x = gen_input()  
    f(x)
```

Human Effort: figure out **input-outputs**, write the **test cases**

Human Effort: figure out the **very bad input**, describe the **bad behavior**

Penetration testing

Property-based testing

Randomly generating  
test inputs, but with  
properties checking outputs

```
for i in 0..10^4:  
    x = gen_input()  
    y = f(x)  
    assert property(x,y)
```

Crafting adversarial inputs  
to trigger vulnerabilities rather  
than getting outputs

```
try:  
    f(very_bad_input)  
    assert false  
except:  
    assert no_bad_behavior
```

Human Effort: write the **input generator**, **harness** and the **property**

# Unit Testing: Precision Shots

- Human effort
  - Figure out expected input/output pairs (multiple)
  - Write a test case for each input/output pairs (multiple)
- Goal
  - Maximize confidence via coverage (approx.: line coverage)
  - Execute precise, human-chosen inputs that exercise specific code paths
  - Error handling is a path too!

Program – Line coverage: 0/5 (0%)

```
def calculate_interest(balance: float, months: int) -> float:  
    rate = 0.02 # 2% per month  
    if months > 12:  
        rate = 0.03 # higher rate for long-term deposits  
    return balance * rate / months
```

Program – Line coverage: 0/5 (0%)

```
def calculate_interest(balance: float, months: int) -> float:  
    rate = 0.02 # 2% per month  
    if months > 12:  
        rate = 0.03 # higher rate for long-term deposits  
    return balance * rate / months
```

Test 1

```
def test_short_term_interest():  
    assert calculate_interest(1000, 6) == pytest.approx(3.33)
```

Program – Line coverage: 3/5 (60%)

```
def calculate_interest(balance: float, months: int) -> float:  
    rate = 0.02 # 2% per month  
    if months > 12:  
        rate = 0.03 # higher rate for long-term deposits  
    return balance * rate / months
```

Test 1

```
def test_short_term_interest():  
    assert calculate_interest(1000, 6) == pytest.approx(3.33)
```

Program – Line coverage: 3/5 (60%)

```
def calculate_interest(balance: float, months: int) -> float:  
    rate = 0.02 # 2% per month  
    if months > 12:  
        rate = 0.03 # higher rate for long-term deposits  
    return balance * rate / months
```

Test 1

```
def test_short_term_interest():  
    assert calculate_interest(1000, 6) == pytest.approx(3.33)
```

Test 2

```
def test_long_term_interest():  
    assert calculate_interest(1000, 24) == pytest.approx(1000 * 0.03 / 24)
```

Program – Line coverage: 5/5 (100%)

```
def calculate_interest(balance: float, months: int) -> float:  
    rate = 0.02 # 2% per month  
    if months > 12:  
        rate = 0.03 # higher rate for long-term deposits  
    return balance * rate / months
```

Test 1

```
def test_short_term_interest():  
    assert calculate_interest(1000, 6) == pytest.approx(3.33)
```

Test 2

```
def test_long_term_interest():  
    assert calculate_interest(1000, 24) == pytest.approx(1000 * 0.03 / 24)
```

Are we done covering everything?

Program – Line coverage: 5/5 (100%)

```
def calculate_interest(balance: float, months: int) -> float:  
    rate = 0.02 # 2% per month  
    if months > 12:  
        rate = 0.03 # higher rate for long-term deposits  
    return balance * rate / months
```

Test 1

```
def test_short_term_interest():  
    assert calculate_interest(1000, 6) == pytest.approx(3.33)
```

Test 2

```
def test_long_term_interest():  
    assert calculate_interest(1000, 24) == pytest.approx(1000 * 0.03 / 24)
```

Are we done covering everything?

Program – Line coverage: 5/5 (100%)

```
def calculate_interest(balance: float, months: int) -> float:  
    rate = 0.02 # 2% per month  
    if months > 12:  
        rate = 0.03 # higher rate for long-term deposits  
    return balance * rate / months
```

No!

Test 1

```
def test_short_term_interest():  
    assert calculate_interest(1000, 6) == pytest.approx(3.33)
```

Test 2

```
def test_long_term_interest():  
    assert calculate_interest(1000, 24) == pytest.approx(1000 * 0.03 / 24)
```

Are we done covering everything?

Program – Line coverage: 5/5 (100%)

```
def calculate_interest(balance: float, months: int) -> float:  
    rate = 0.02 # 2% per month  
    if months > 12:  
        rate = 0.03 # higher rate for long-term deposits  
    return balance * rate / months
```

No!

Test 1

```
def test_short_term_interest():  
    assert calculate_interest(100)
```

Test 2

```
def test_long_term_interest():  
    assert calculate_interest(1000, 24) == pytest.approx(1000 * 0.03 / 24)
```

## CWE-369: Divide By Zero

**Weakness ID:** 369

**Vulnerability Mapping:** ALLOWED

**Abstraction:** Base

Are we done covering everything?

Program – Line coverage: 5/5 (100%)

```
def calculate_interest(balance: float, months: int) -> float:  
    rate = 0.02 # 2% per month  
    if months > 12:  
        rate = 0.03 # higher rate for long-term deposits  
    return balance * rate / months
```

No!

Test 1

```
def test_short_term_interest():  
    assert calculate_interest(1000, 6) == pytest.approx(3.33)
```

Test 2

```
def test_long_term_interest():  
    assert calculate_interest(1000, 24) == pytest.approx(1000 * 0.03 / 24)
```

Test 3

```
def test_zero_months_divide_by_zero(): # covers implicit error path  
    with pytest.raises(ZeroDivisionError):  
        calculate_interest(1000, 0)
```

# **Automated Unit Test Improvement using Large Language Models at Meta**

Nadia Alshahwan\*

Jubin Chheda

Anastasia Finegenova

Beliz Gokkaya

Mark Harman

Inna Harper

Alexandru Marginean

Shubho Sengupta

Eddy Wang

Meta Platforms Inc.,

Menlo Park, California, USA

## Automated Unit Test Improvement using Large Language Models at Meta

Nadia Alshahwan\*

Jubin Chheda

Anastasia Finegenova

Beliz Gokkaya

Mark Harman

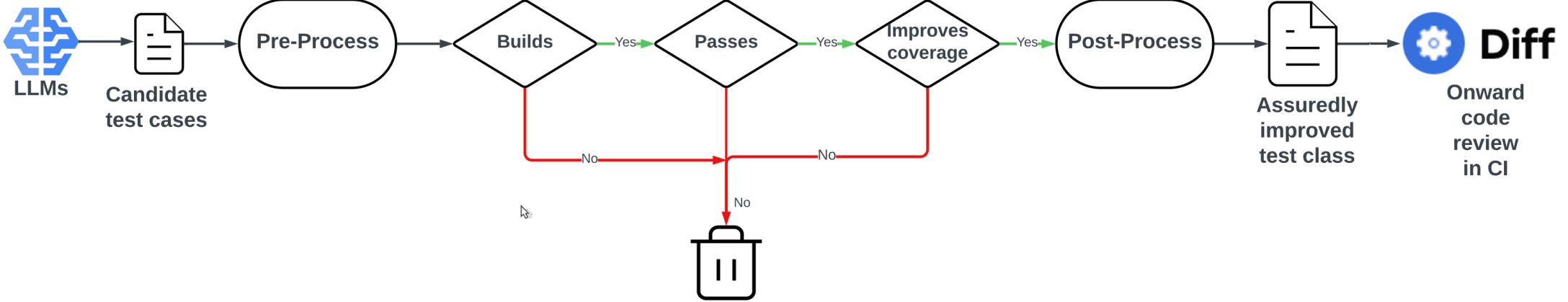
Inna Harper

Alexandru Marginean

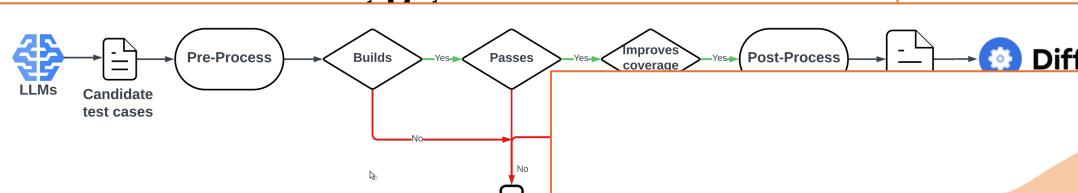
Shubho Sengupta

Eddy Wang

Meta Platforms Inc.



## Automated Unit Test Improvement using Large Language Models



Ludy Wang  
Meta Platforms Inc.,  
Menlo Park, California, USA

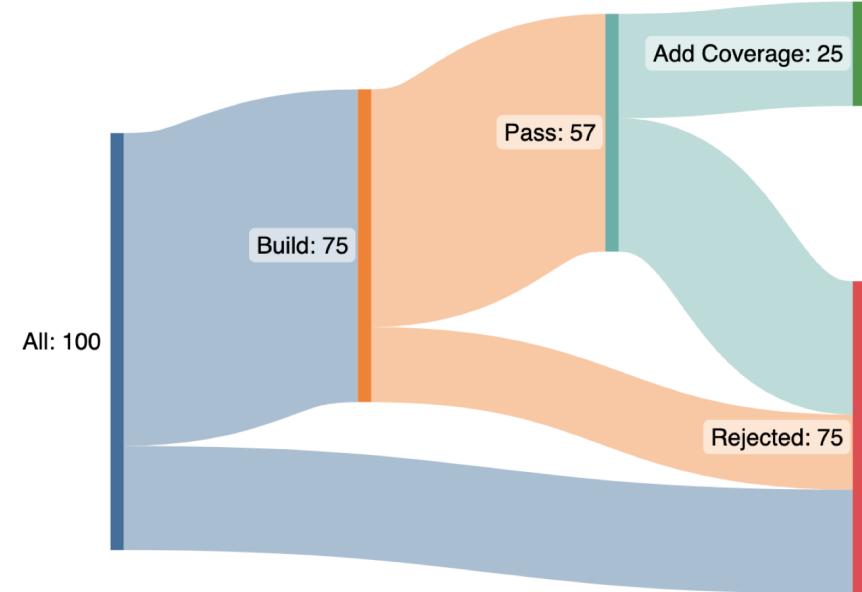


Figure 2: Sankey diagram showing the filtration process outcomes (as percentages of all test cases) from the Experimental Study on Instagram components for Reels and Stories products, using the four prompt strategies from Table 2 and the two language models, LLM1 and LLM2.

## Automated Unit Test

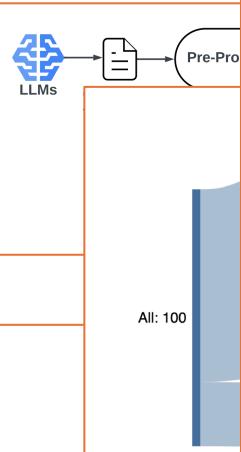


Figure 2: Success rates of generated test cases (as percentage) from the Study on Instantaneous Test Generation using Large Language Models, using the two language models.

### 3.4.1 Deployment at the Facebook App Test-a-thon December 2023.

In this deployment, we had sufficient confidence to automatically submit recommendations from TestGen-LLM to engineers. There was no engineer pre-training process, no specific test-a-thon expectations, and no additional context provided to the engineers. This gave us a realistic assessment of the engineers' response to LLM-generated test recommendations provided 'out of the box'. Overall, over 50% of the diffs submitted were accepted by developers, a figure which rises to almost 70% of those which received a review by developers. Specifically, of the 280 diffs generated by TestGen-LLM:

- 144 were accepted by the engineer reviewing them.
- 64 were rejected and/or abandoned.
- 61 did not receive a review.
- 11 were withdrawn.

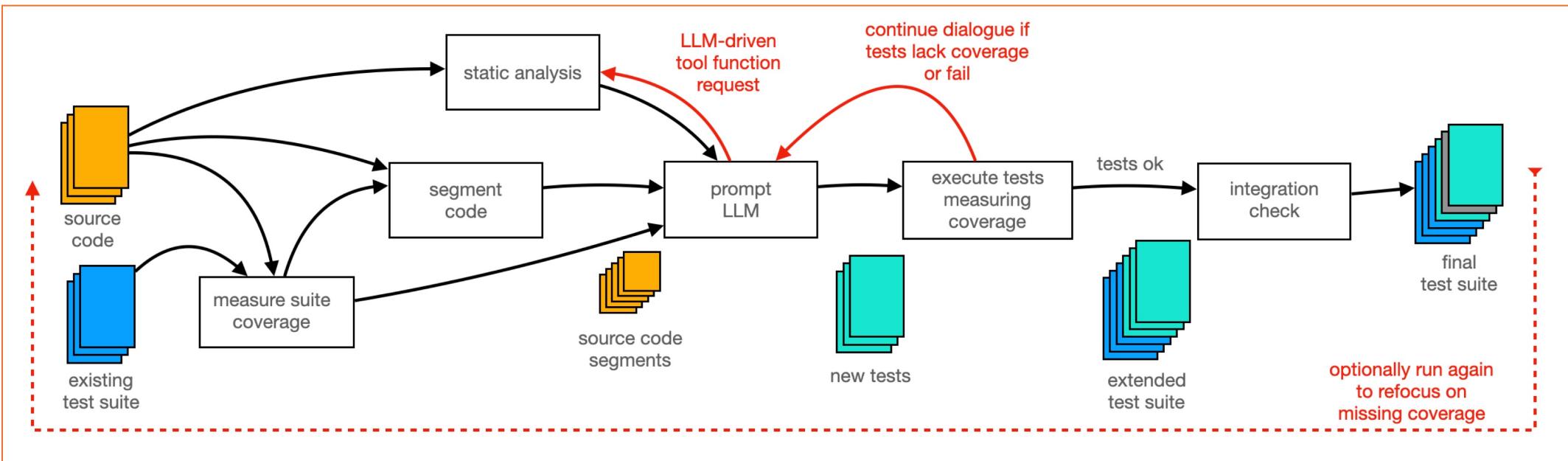
## **COVERUP: Coverage-Guided LLM-Based Test Generation**

JUAN ALTMAYER PIZZORNO, University of Massachusetts Amherst, United States

EMERY D. BERGER\*, University of Massachusetts Amherst / Amazon, United States

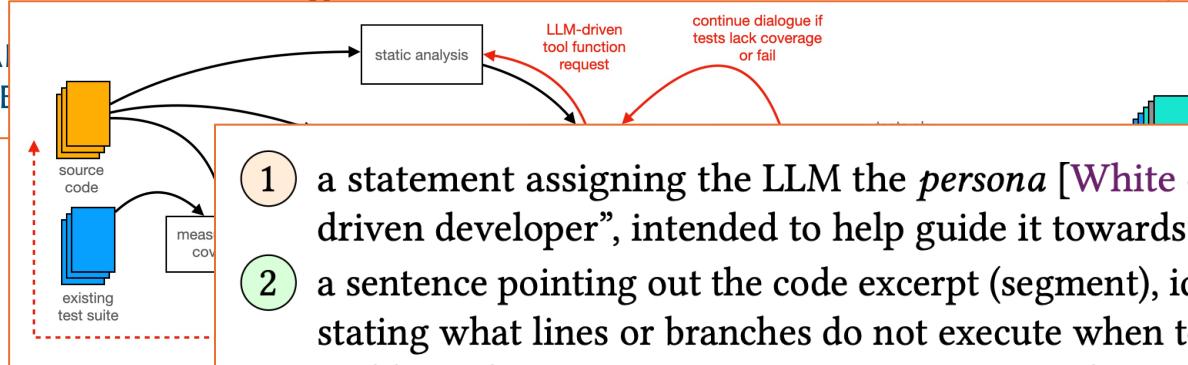
# COVERUP: Coverage-Guided LLM-Based Test Generation

JUAN ALTMAYER PIZZORNO, University of Massachusetts Amherst, United States  
EMERY D. BERGER\*, University of Massachusetts Amherst / Amazon, United States



## COVERUP: Coverage-Guided LLM-Based Test Generation

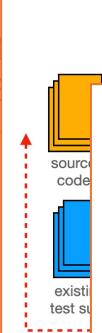
JUAI  
EME



- 1 a statement assigning the LLM the *persona* [White et al. 2023] of an “expert Python test-driven developer”, intended to help guide it towards high quality tests;
- 2 a sentence pointing out the code excerpt (segment), identifying what file it comes from, and stating what lines or branches do not execute when tested. The portion specifying the lines and branches missing coverage is compressed using line ranges, simplifying the prompt and reducing token usage.
- 3 a request for pytest test functions and an encouragement for the LLM to use the provided tool function;
- 4 a series of other requests, such as “include assertions” and “avoid state pollution”, to steer the result towards usable tests;
- 5 a request that the response only include the new Python tests to facilitate its extraction from the response and to reduce token usage; and
- 6 the code segment, prefixed by generated import statements and tagging the lines lacking (line or branch) coverage with their numbers.

## COVERUP: Coverage-Guided LLM-Based Test Generation

JUAI  
EME



- 1 a statement assigning the LLM the *persona* [White et al. 2023] of an “expert Python test-driven developer”, intended to help guide it towards high quality tests;
- 2 a sentence pointing out the code excerpt (segment), identifying what file it comes from, and ~~stating what lines or branches do not execute when tested. The portion specifying the lines~~

(1)

This test still lacks coverage: line 615 and branches 603->exit, 610->608, 618->exit do not execute.

(2)

Modify it to correct that; respond only with the complete Python code in backticks.  
Use the `get_info` function as necessary.

Fig. 6. **Example of a coverage follow-up prompt:** COVERUP indicates to the LLM that a line and some branches still weren't covered (1), asking that it correct the test (2).

## COVERUP: Coverage-Guided LLM-Based Test Generation

JUAI  
EME

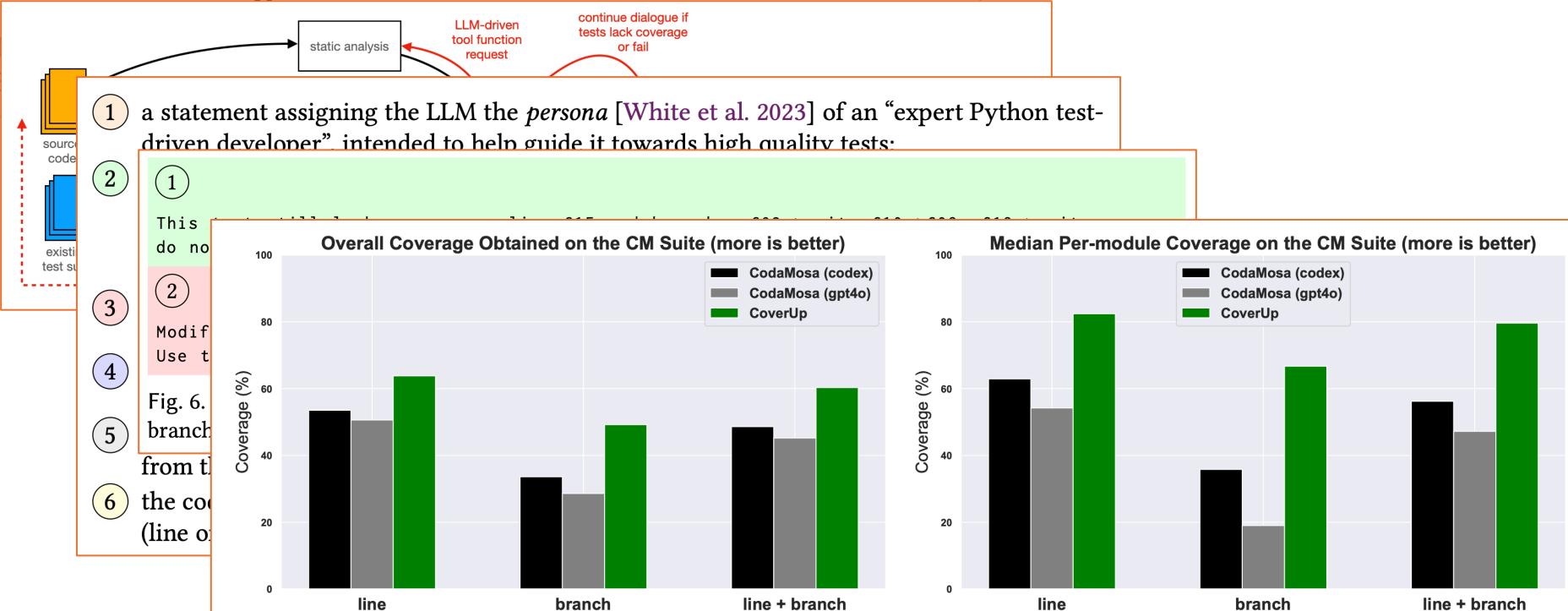


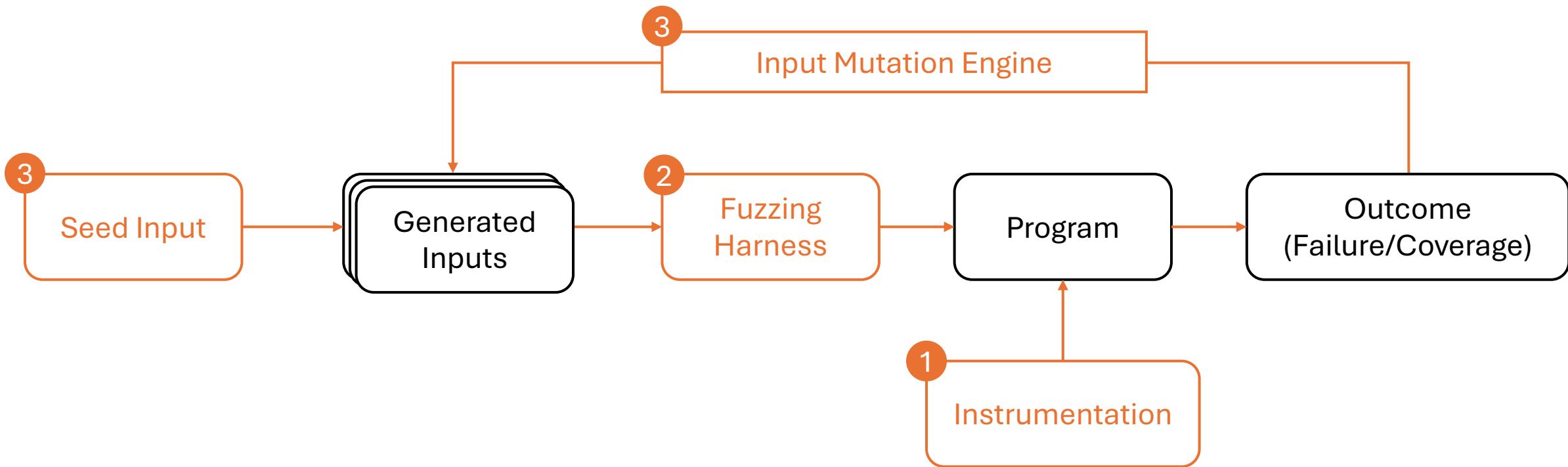
Fig. 12. [RQ1] **COVERUP yields higher overall and median per-module coverage:** Across the board, COVERUP yields higher coverage than CODAMOSA, whether measured over the entire suite or on a module-by-module basis.

# Fuzzing: Chaos with Purpose

```
for i in 0..10^4:  
    x = gen_input()  
    f(x)
```

- Goal
  - Explore huge input spaces automatically to trigger failures (crashes, exceptions, timeouts, sanitizer hits)
  - Unlike unit tests, fuzzing does not know the correct output. It detects incorrect behavior by symptoms
- Method
  - Massively and randomly generate inputs to programs

# Fuzzing: The Loop



```
int process_input(const uint8_t *data, size_t size) {
    if (size < 2) return 0; // need at least 2 bytes: cmd + length
    uint8_t cmd = data[0];
    uint8_t n = data[1]; // number of bytes to copy

    char buf[16]; // fixed-size buffer
    int secret = 0x41414141; // adjacent variable

    // If cmd == 0x42 we write to buffer
    if (cmd == 0x42) {
        // Write to buffer
        buf[n] = 0x00;
    }

    // Return secret so program behavior may subtly change if corrupted
    return secret;
}
```

```

int process_input(const uint8_t *data, size_t size) {
    if (size < 2) return 0; // need at least 2 bytes: cmd + length
    uint8_t cmd = data[0];
    uint8_t n = data[1]; // number of bytes to copy

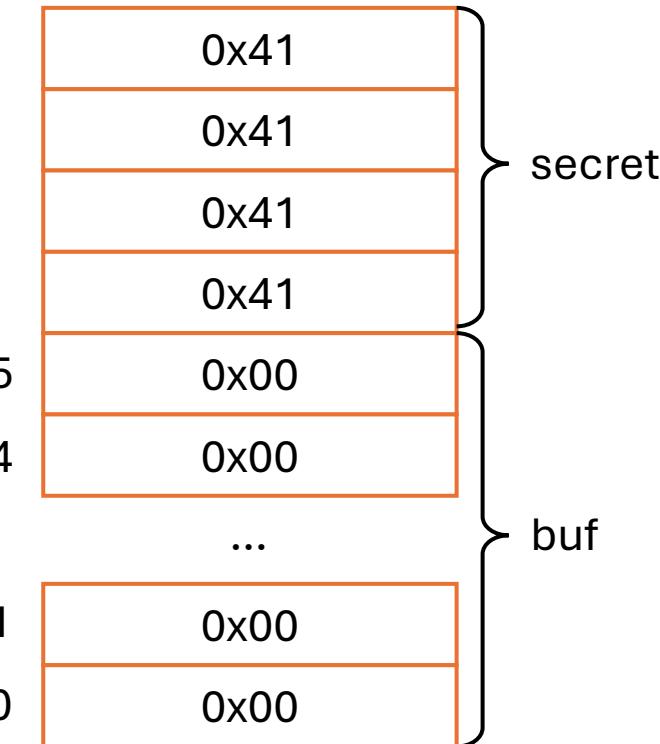
    char buf[16]; // fixed-size buffer
    int secret = 0x41414141; // adjacent variable

    // If cmd == 0x42 we write to buffer
    if (cmd == 0x42) {
        // Intentionally naive: no bounds check on n relative to buf size
        buf[n] = 0x00;
    }

    // Return secret so program behavior may subtly change if corrupted
    return secret;
}

```

Memory (Stack in bytes)



```

int process_input(const uint8_t *data, size_t size) {
    if (size < 2) return 0; // need at least 2 bytes: cmd + length
    uint8_t cmd = data[0];
    uint8_t n = data[1]; // number of bytes to copy

    char buf[16]; // fixed-size buffer
    int secret = 0x41414141; // adjacent variable

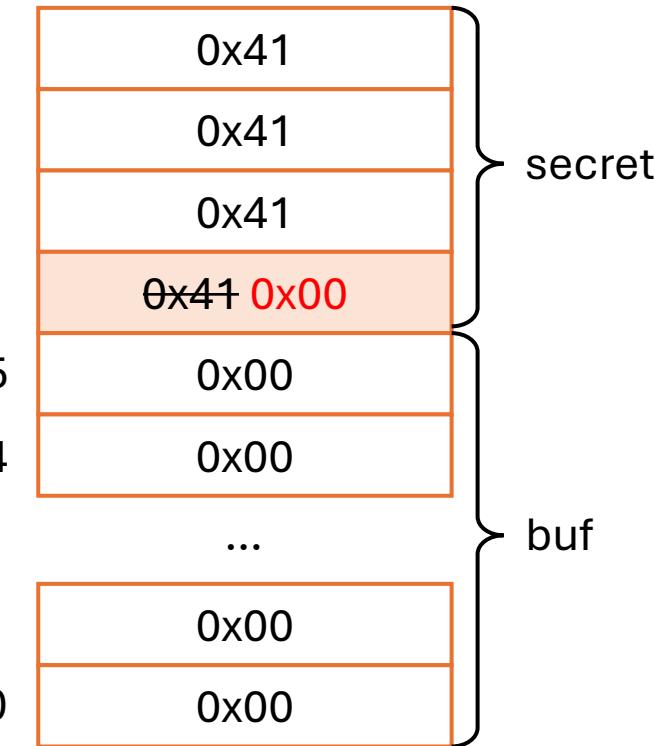
    // If cmd == 0x42 we write to buffer
    if (cmd == 0x42) {
        // Intentionally naive: no bounds check on n relative to buf size
        buf[n] = 0x00;
    }

    // Return secret so program behavior may subtly change if corrupted
    return secret;
}

```

process\_input({0x42, 16}, 2)

Memory (Stack in bytes)



```

int process_input(const uint8_t *data, size_t size) {
    if (size < 2) return 0; // need at least 2 bytes: cmd + length
    uint8_t cmd = data[0];
    uint8_t n = data[1]; // number of bytes to copy

    char buf[16]; // fixed-size buffer
    int secret = 0x41414141; // adjacent variable

    // If cmd == 0x42 we write to buffer
    if (cmd == 0x42) {
        // Intentionally naive: no bounds check on n relative to buf size
        buf[n] = 0x00;
    }

    // Return secret so program behavior may subtly change if corrupted
    return secret;
}

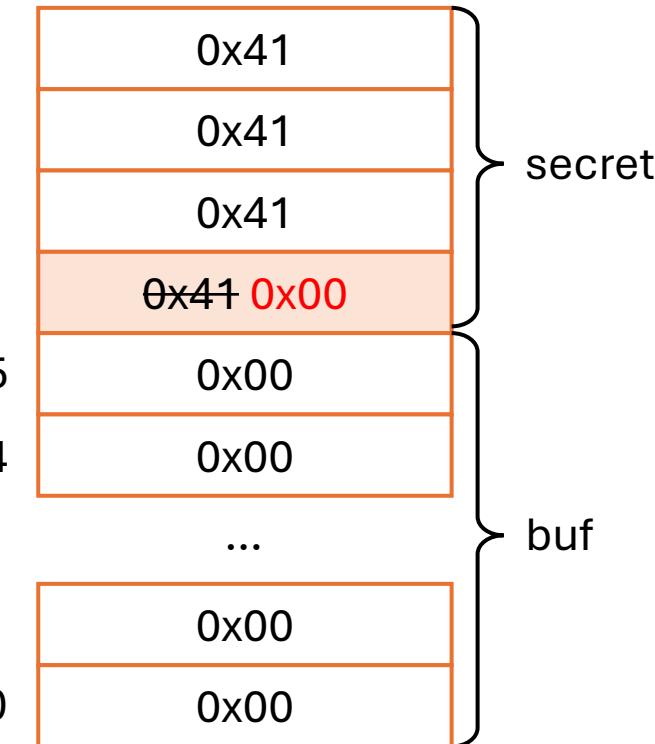
```

process\_input({0x42, 16}, 2)

#### Issue:

- There is an **out-of-bounds buffer** write!
- But the program will not crash!
- How do we still detect the bug?

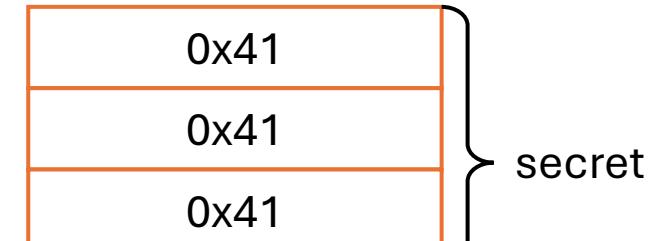
Memory (Stack in bytes)



```
int process_input(const uint8_t *data, size_t size) {
    if (size < 2) return 0; // need at least 2 bytes: cmd + length
    uint8_t cmd = data[0];
    uint8_t n = data[1]; // number of bytes to copy

    char buf[16]; // fixed-size buffer
```

Memory (Stack in bytes)



## 2024 CWE Top 25 Most Dangerous Software Weaknesses

[Top 25 Home](#)

Share via:

[View in table format](#)

[Key Insights](#)

[Methodology](#)

1

Improper Neutralization of Input During Web Page Generation ('Cross-site Scripting')  
[CWE-79](#) | CVEs in KEV: 3 | Rank Last Year: 2 (up 1)

2

Out-of-bounds Write  
[CWE-787](#) | CVEs in KEV: 18 | Rank Last Year: 1 (down 1)

Issue:

- There is an **out-of-bounds buffer** write!
- But the program will not crash!
- How do we still detect the bug?

```

int process_input(const uint8_t *data, size_t size) {
    if (size < 2) return 0; // need at least 2 bytes: cmd + length
    uint8_t cmd = data[0];
    uint8_t n = data[1]; // number of bytes to copy

    char buf[16]; // fixed-size buffer
    int secret = 0x41414141; // adjacent variable

    // If cmd == 0x42 we write to buffer
    if (cmd == 0x42) {
        // Intentionally naive: no bounds check on n relative to buf size
        buf[n] = 0x00;
    }

    // Return secret so program behavior may subtly change if corrupted
    return secret;
}

```

process\_input({0x42, 16}, 2)

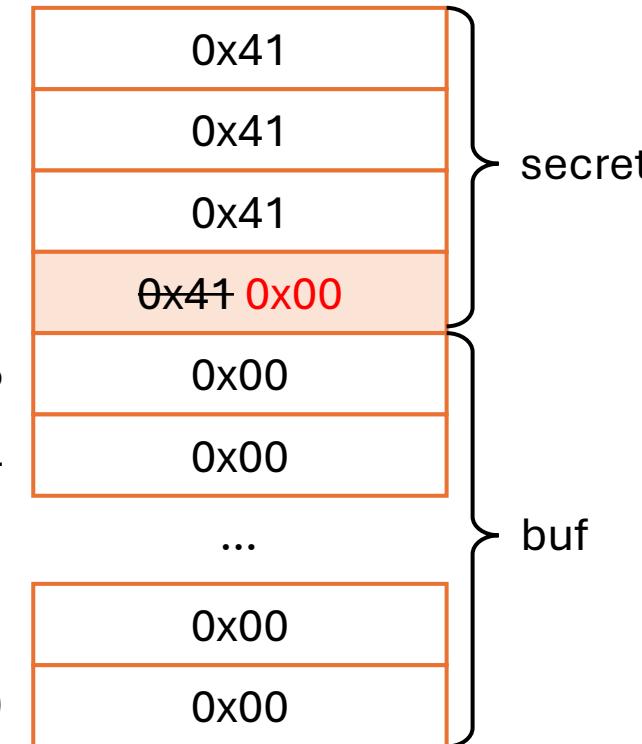
#### Issue:

- There is an **out-of-bounds buffer** write!
- But the program will not crash!
- How do we still detect the bug?

#### Solution:

- Add **Sanitizer!**

Memory (Stack in bytes)



```

int process_input(const uint8_t *data, size_t size) {
    if (size < 2) return 0; // need at least 2 bytes: cmd + length
    uint8_t cmd = data[0];
    uint8_t n = data[1]; // number of bytes to copy

    char buf[16]; // fixed-size buffer
    int secret = 0x41414141; // adjacent variable

    // If cmd == 0x42 we write to buffer
    if (cmd == 0x42) {
        if (n >= 16) SANITIZER_ERROR("out-of-bounds write")
        buf[n] = 0x00;
    }

    // Return secret so program behavior may subtly change if corrupted
    return secret;
}

```

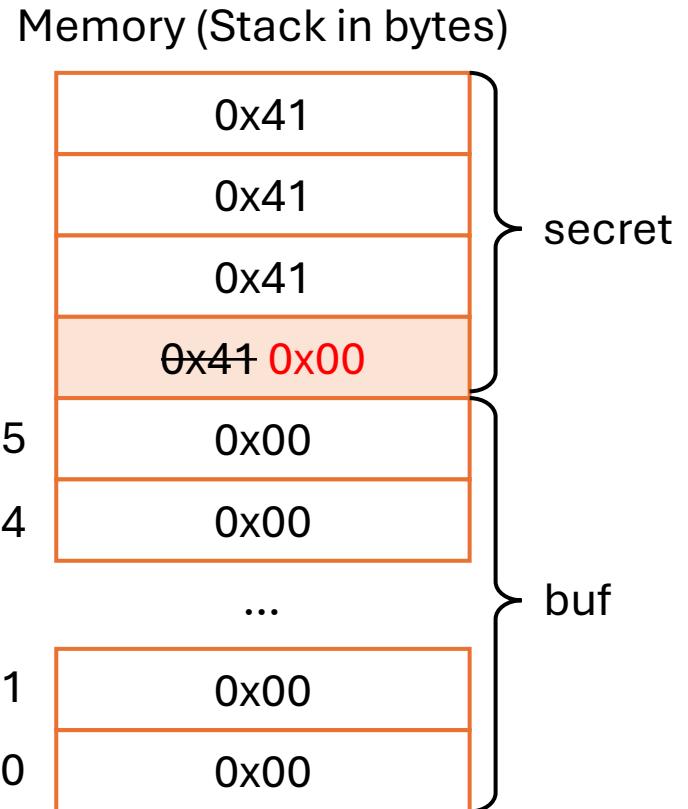
process\_input({0x42, 16}, 2)

#### Issue:

- There is an **out-of-bounds buffer** write!
- But the program will not crash!
- How do we still detect the bug?

#### Solution:

- Add **Sanitizer!**



# Program Instrumentation

 $P$ 

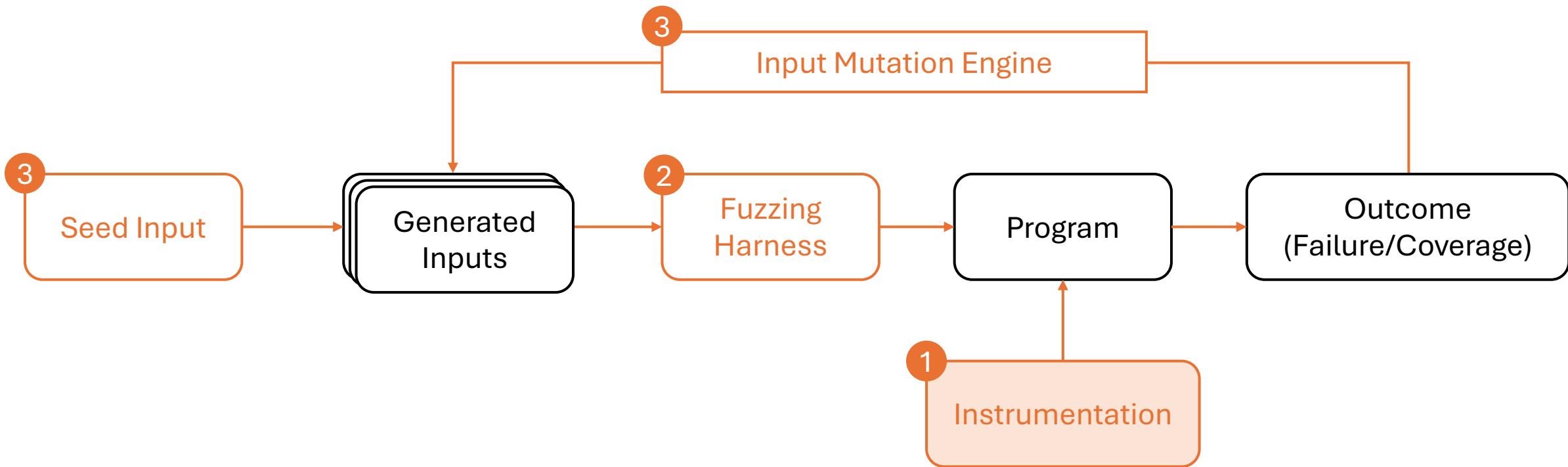
Original Program

 $P'$ 

Instrumented Program

$P$  and  $P'$  are equivalent under normal behavior  
 $P'$  can expose more failure modes than  $P$

# Fuzzing: The Loop



# Fuzzing Harness

- For unit tests:
  - We come up with a single input-output pair
  - We write a test case that runs the program on input, compares the output
- For fuzzing:
  - We have automatically generated input
  - We write a test case that runs the program on input // ← HARNESS

# Fuzzing Harness

- For unit tests:
  - We come up with a single input-output pair
  - We write a test case that runs the program on input, compares the output
- For fuzzing:
  - We have automatically generated input
  - We write a test case that runs the program on input // ← HARNESS

We have to do this  
for a lot of times!

# Fuzzing Harness

- For unit tests:
  - We come up with a single input-output pair
  - We write a test case that runs the program on input, compares the output
- For fuzzing:
  - We have automatically generated input
  - We write a test case that runs the program on input // ← HARNESS

We have to do this  
for a lot of times!

Can we just do this one  
time or only a few times?

 pnggroup / libpng

Type  to search

 Watch 53   Fork 721   Star 1.5k 

 Code  Issues 144  Pull requests 20  Discussions  Actions  Projects  Wiki  Security  Insights

 libpng Public

 libpng16  11 Branches  1625 Tags  Go to file  Add file  Code

This branch is [11 commits ahead of](#), [99 commits behind](#) libpng18.

 stoeckmann and ctruta api! Remove the experimental (and incomplete)...	  bd41aa6 · 3 weeks ago	 4,459 Commits
.github/workflows	ci: Add GitHub Actions for verifying libpng on Linux, mac...	last month
arm	chore: Clean up the error directives	6 months ago
ci	ci: Run autogen.sh without --maintainer in ci_verify_con...	last month
contrib	refactor: Delete conditional compilation for libpng 1.6.0 or...	last month
intel	[libpng16] chore: Clean up the leading blank lines from all...	last year

 About

LIBPNG: Portable Network Graphics support, official libpng repository

 [libpng.sf.net](#)

 Readme

 View license

 Activity

 Custom properties

 1.5k stars

 53 watching

 721 forks

Report repository

<https://github.com/pnggroup/libpng>

This branch is 11 commits ahead of, 99 commits behind libpng18.

Name	Last commit message
..	
Dockerfile	Revert "oss-fuzz: Update the README file, the Docker file and the bui...
README.txt	Revert "oss-fuzz: Update the README file, the Docker file and the bui...
build.sh	Revert "oss-fuzz: Update the README file, the Docker file and the bui...
libpng_read_fuzzer.cc	[libpng16] chore: Clean up the leading blank lines from all source files
libpng_read_fuzzer.options	Revert "oss-fuzz: Transfer to an external repo and remove from this r..."
png.dict	Revert "oss-fuzz: Transfer to an external repo and remove from this r..."

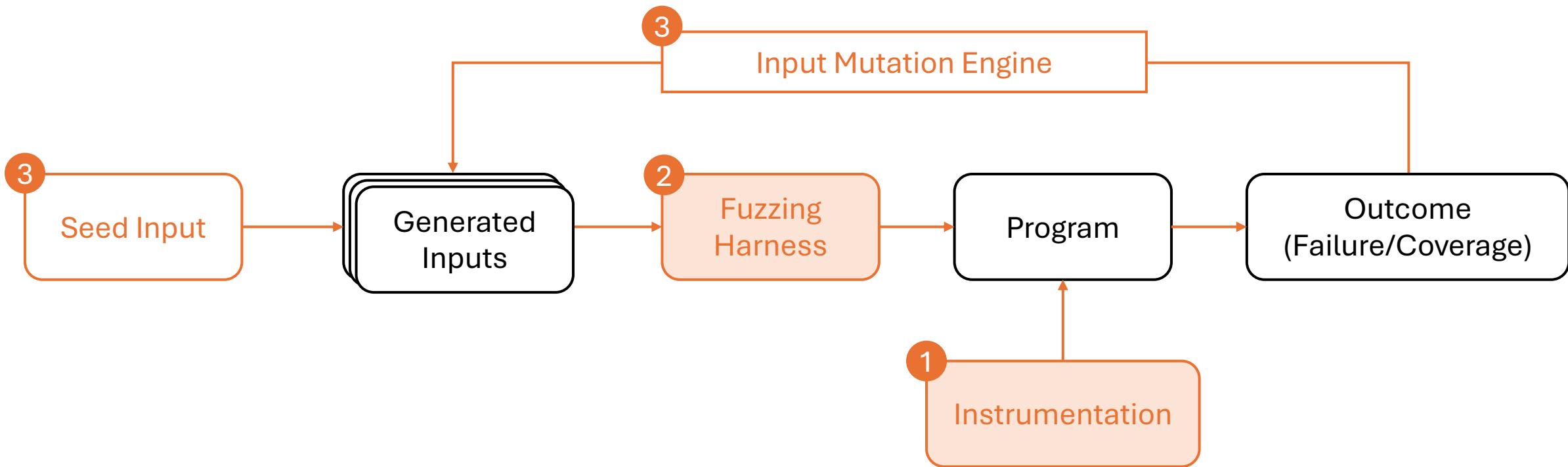
The screenshot shows the GitHub repository for libpng. It displays two branches: 'libpng16' (selected) and 'libpng18'. The 'libpng16' branch has 144 issues, 20 pull requests, and 11 commits ahead of 'libpng18'. The repository contains files like Dockerfile, README.txt, build.sh, libpng\_read\_fuzzer.cc, libpng\_read\_fuzzer.options, and png.dict.

```
97     // Entry point for LibFuzzer.
98     // Roughly follows the libpng book example:
99     // http://www.libpng.org/pub/png/book/chapter13.html
100    extern "C" int LLVMFuzzerTestOneInput(const uint8_t* data, size_t size) {
101        if (size < kPngHeaderSize) {
102            return 0;
103        }
104
105        std::vector<unsigned char> v(data, data + size);
106        if (png_sig_cmp(v.data(), 0, kPngHeaderSize)) {
107            // not a PNG.
108            return 0;
109        }
110
111        PngObjectHandler png_handler;
112        png_handler.png_ptr = nullptr;
113        png_handler.row_ptr = nullptr;
114        png_handler.info_ptr = nullptr;
115        png_handler.end_info_ptr = nullptr;
116
117        png_handler.png_ptr = png_create_read_struct(
118            PNG_LIBPNG_VER_STRING, nullptr, nullptr, nullptr);
119        if (!png_handler.png_ptr) {
120            return 0;
121        }
122
123        png_handler.info_ptr = png_create_info_struct(png_handler.png_ptr);
124        if (!png_handler.info_ptr) {
125            PNG_CLEANUP
126            return 0;
127        }
128
129        png_handler.end_info_ptr = png_create_info_struct(png_handler.png_ptr);
130        if (!png_handler.end_info_ptr) {
131            PNG_CLEANUP
132            return 0;
133        }
```

# Fuzzing Harness: Key Message

- A single wide-capturing test case that takes a single buffer as input
  - Reject bad inputs that we do not care about in the beginning
  - Simulate complete usage of a library or a package
  - Trigger sanitization errors or capture failures and exceptions
- In practice, a library may be accompanied multiple harnesses
  - But not to the level of the amount of unit test cases

# Fuzzing: The Loop



# Seed Input and Mutation Engine

- Seed Input
  - Typically some valid inputs to a program
  - E.g., an HTML for a browser, a PNG for libpng, a JSON for a JSON parser
- Mutation Engine
  - With a single input, we may get **feedback** from the execution:
    - E.g., coverage, error, command line outputs, etc.
  - We may use the coverage to guide the mutation engine to **change** the seed input to something else
    - E.g., JSON “[1, 2, 3]” → “[1, 2, ]”
  - **Hopefully**, the mutation leads us to **a better coverage**

# Infinite Monkey Theorem

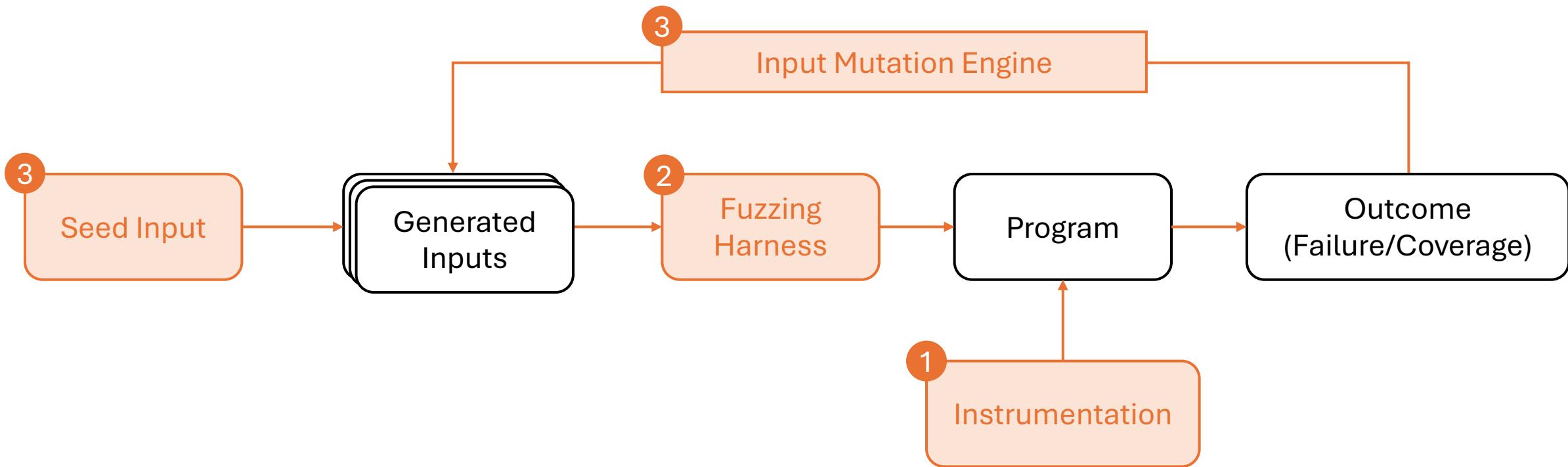
*“a monkey hitting keys independently  
and at random on a typewriter keyboard  
for an infinite amount of time will almost  
surely type any given text, including the  
complete works of William  
Shakespeare.”*



# Seed Input and Mutation Engine

- Seed Input
  - Typically some valid inputs to a program
  - E.g., an HTML for a browser, a PNG for libpng, a JSON for a JSON parser
- Mutation Engine
  - With a single input, we may get **feedback** from the execution:
    - E.g., coverage, error, command line outputs, etc.
  - We may use the coverage to guide the mutation engine to **change** the seed input to something else
    - E.g., JSON “[1, 2, 3]” → “[1, 2, ]”
  - **Hopefully**, the mutation leads us to **a better coverage**

# Fuzzing: The Loop



# Fuzzing in the Wild: AFL

**google/AFL**

american fuzzy lop - a security-oriented fuzzer

“American Fuzzing Loop”

24 Contributors    54 Issues    4k Stars    663 Forks



# AFL++: Combining Incremental Steps of Fuzzing Research

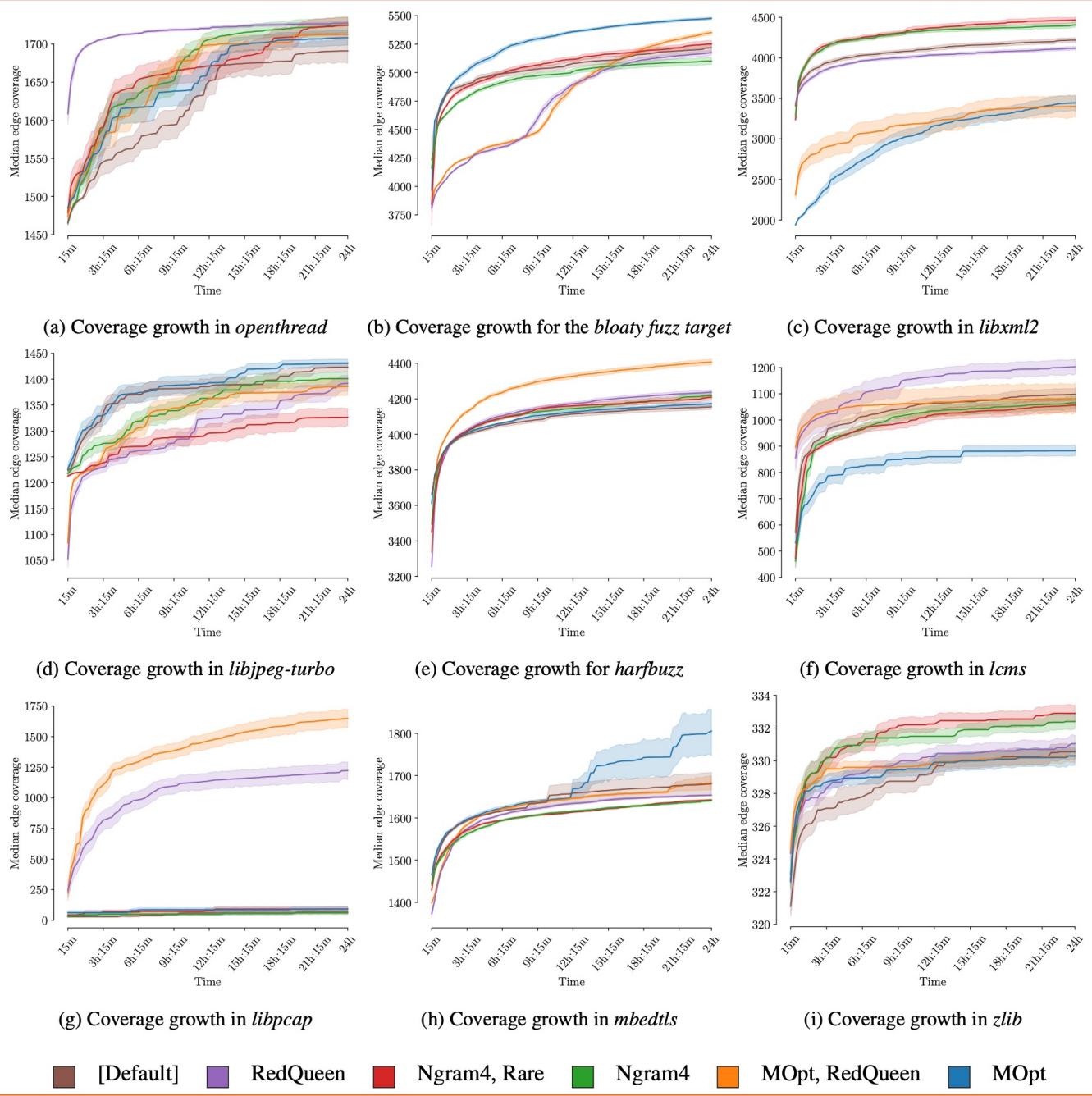
Andrea Fioraldi<sup>†</sup>, Dominik Maier<sup>‡</sup>, Heiko Eißfeldt, Marc Heuse<sup>§</sup>

*{andrea, dominik, heiko, marc}@aflplus.plus*

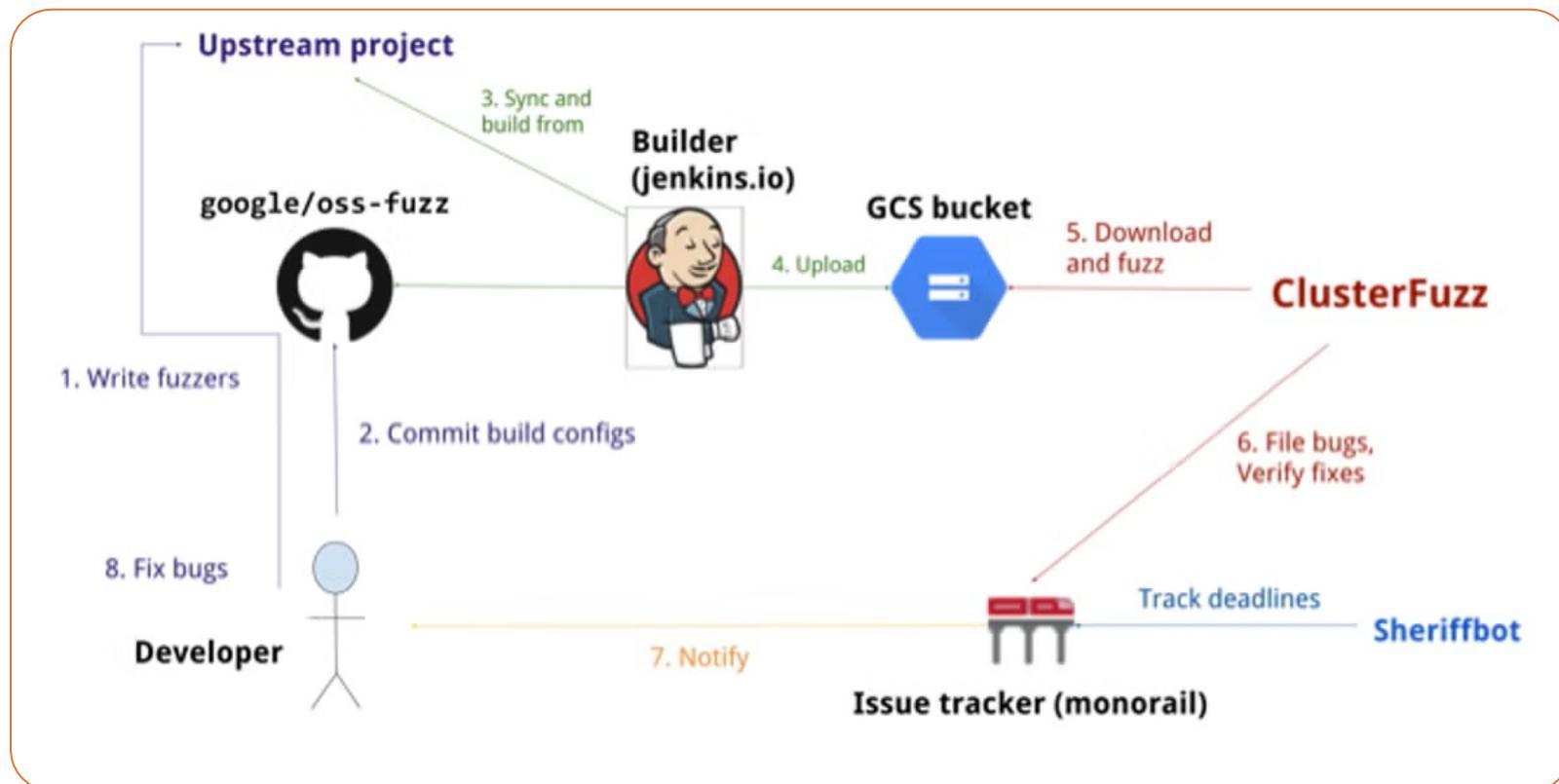
*†Sapienza University of Rome, ‡TU Berlin, §The Hacker’s Choice*

# AFL++: Combining Incremental Coverage

Andrea Fioraldi<sup>†</sup>, Dominik Kullmann<sup>‡</sup>  
*{andrea, dominik}@csail.mit.edu*  
†Sapienza University of Rome, Italy



# Fuzzing in the Wild: OSS-Fuzz



# Fuzzing in the Wild: OSS-Fuzz

## Coverage Report

View results by: [Directories](#) | [Files](#)

PATH	LINE COVERAGE	FUNCTION COVERAGE	REGION COVERAGE
<a href="#">autofit/</a>	89.65% (6801/7586)	88.64% (78/88)	86.09% (3280/3810)
<a href="#">base/</a>	73.47% (7650/10413)	77.81% (249/320)	62.89% (5010/7966)
<a href="#">bdf/</a>	94.60% (2503/2646)	97.30% (36/37)	92.54% (2133/2305)
<a href="#">cff/</a>	91.68% (5048/5506)	92.08% (93/101)	89.75% (2880/3209)
<a href="#">cid/</a>	96.03% (1717/1788)	93.94% (31/33)	94.75% (1100/1161)
<a href="#">gzip/</a>	95.54% (1435/1502)	94.12% (32/34)	88.73% (1299/1464)
<a href="#">lzw/</a>	79.80% (482/604)	83.33% (15/18)	76.90% (263/342)
<a href="#">pcf/</a>	94.46% (1722/1823)	83.87% (26/31)	93.17% (1445/1551)
<a href="#">pfr/</a>	91.69% (2603/2839)	91.49% (43/47)	92.17% (2049/2223)
<a href="#">psaux/</a>	91.19% (11504/12615)	96.12% (198/206)	88.69% (5990/6754)
<a href="#">pshinter/</a>	95.20% (2957/3106)	92.41% (73/79)	91.25% (1356/1486)
<a href="#">psnames/</a>	99.36% (466/469)	100.00% (11/11)	98.70% (227/230)
<a href="#">raster/</a>	96.07% (2002/2084)	90.48% (38/42)	91.34% (1539/1685)
<a href="#">sfnt/</a>	89.03% (8089/9086)	96.00% (144/150)	87.35% (6187/7083)
<a href="#">smooth/</a>	90.62% (899/992)	82.14% (23/28)	85.52% (685/801)
<a href="#">truetype/</a>	95.29% (12299/12907)	99.63% (269/270)	89.28% (7409/8299)
<a href="#">type1/</a>	93.10% (4154/4462)	93.94% (62/66)	90.94% (2811/3091)
<a href="#">type42/</a>	95.79% (1662/1735)	90.62% (29/32)	94.50% (1099/1163)
<a href="#">winfonts/</a>	96.16% (826/859)	100.00% (13/13)	94.33% (632/670)
<b>TOTALS</b>	<b>90.12% (74819/83022)</b>	<b>91.10% (1463/1606)</b>	<b>85.71% (47394/55293)</b>

```
59      /* multiply a given value by the CORDIC shrink factor */
60      static FT_Fixed
61      ft_trig_downscale( FT_Fixed val )
62      {
63          FT_Int s = 1;
64
65          if ( val < 0 )
66          {
67              val = -val;
68              s = -1;
69          }
70
71
72          /* 0x40000000 comes from regression analysis between true */
73          /* and CORDIC hypotenuse, so it minimizes the error */
74          val = (FT_Fixed)(
75              ( (FT_UInt64)val * FT_TRIG_SCALE + 0x40000000UL ) >> 32 );
76
77          return s < 0 ? -val : val;
78      }
```

# Fuzzing in the Wild: Fuzzing Introspection

**Open Source Fuzzing Introspection**

Project Overview   Function Database   Target Oracle   Indexing status   API   Fuzz Introspector   OSS-Fuzz   About

## Fuzzing Introspection of OSS-Fuzz projects

This page shows stats and data on open source fuzzing of the projects integrated into OSS-Fuzz. The analysis is generated by Fuzz Introspector, which is our tool for analysing the quality of fuzzing for an open source project. The goal is to make the status transparent and useful for developers and researchers to identify if the code they use is properly analysed.

[Projects overview](#)   [Search database](#)

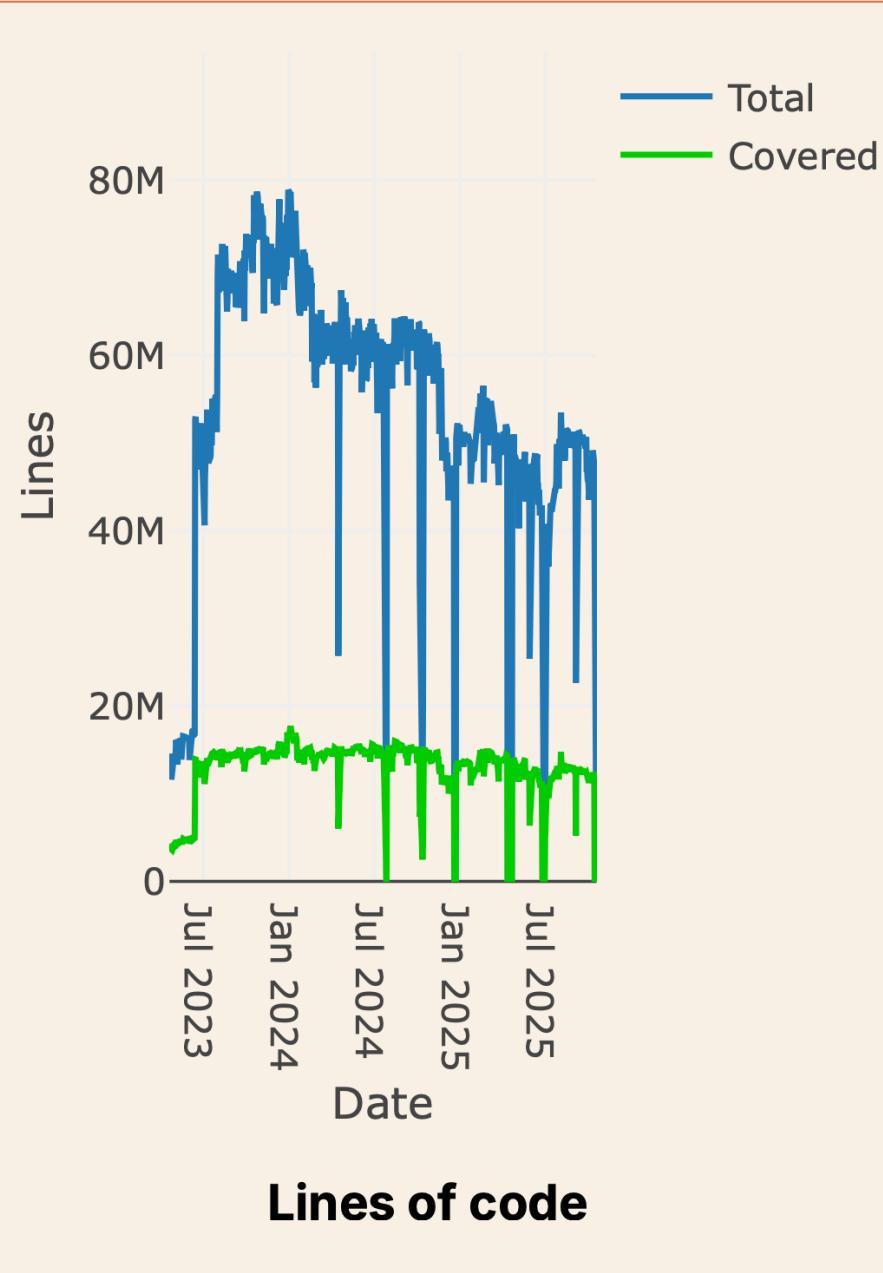
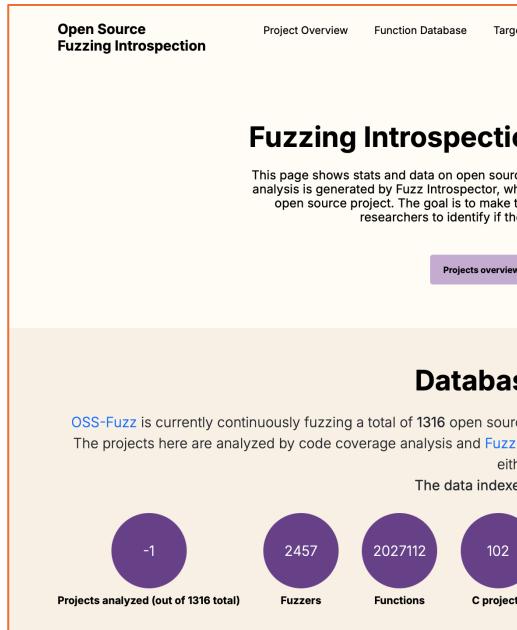
## Database overview

OSS-Fuzz is currently continuously fuzzing a total of 1316 open source projects. The database on this page shows details about a subset of these. The projects here are analyzed by code coverage analysis and Fuzz Introspector tool, and we only show data where there is a successful build of either of these.

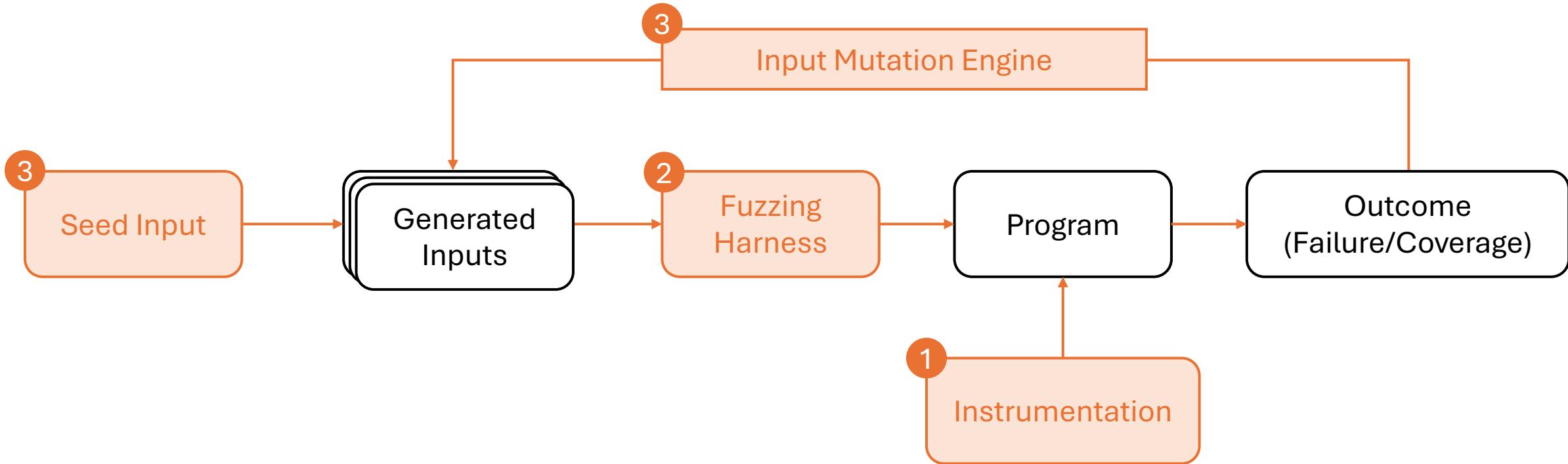
The data indexed on this site comprises

 -1	 2457	 2027112	 102	 262	 221	 97	 0	 0
Projects analyzed (out of 1316 total)	Fuzzers	Functions	C projects	C++ projects	Python projects	Java projects	Go projects	Rust projects

# Fuzzing in the wild: OSS-Fuzz inspection



# How Can LLM Help?



# How Can LLM Help?

- LLM for **instrumentation**:
  - Capture more vulnerability patterns than traditional sanitizers
  - E.g., Path Traversal, Reading from external logs, etc.
- LLM for **fuzzing harness generation**:
  - Capture more library usage pattern
  - E.g., For libpng, simulate browser usage, PDF reader usage, etc.
- LLM for **seed input generation & mutation**:
  - Similar to unit test generation, but can accept more types of feedback
  - E.g., Looking at uncovered lines, mutate input towards those lines

# How Effective Are They? Exploring Large Language Model Based Fuzz Driver Generation

Cen Zhang

Nanyang Technological University  
Singapore

Yaowen Zheng\*

Nanyang Technological University  
Singapore

Mingqiang Bai

IIE, CAS; Sch of Cyber Security, UCAS  
Beijing, China

Yeting Li

IIE, CAS; Sch of Cyber Security, UCAS  
Beijing, China

Wei Ma

Nanyang Technological University  
Singapore

Xiaofei Xie

Singapore Management University  
Singapore

Yuekang Li

The University of New South Wales  
Sydney, Australia

Limin Sun

IIE, CAS; Sch of Cyber Security, UCAS  
Beijing, China

Yang Liu

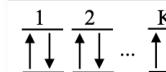
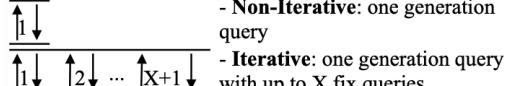
Nanyang Technological University  
Singapore

# How Effective Are They? Exploring Large Language Model Based

Cen Z  
Nanyang Technolo  
Singa

Yetin  
IIE, CAS; Sch of Cyb  
Beijing,

Yueka  
The University of N  
Sydney, A

Legend:	
	Basic API info only
<b>Prompt Strategy Design</b>	
<b>Design I - [Query With Different Types of API Info]</b>	
- <b>Basic Info</b> : Precisely specified & generally accessible, e.g., header file and API declaration	
- <b>Extended Info</b> : Not guaranteed in quality and availability, e.g., API documentation and usage snippets	
<b>Design II - [Query Repeatedly]</b>	
	- Repeat whole query process K times for K independent results - Suffix "-K" in name
<b>Design III - [Query Iteratively]</b>	
	- Non-Iterative: one generation query - Iterative: one generation query with up to X fix queries
<b>Acronym</b>	<b>Prompts (Generation + Fix)</b>
NAIVE-K	①
BACTX-K	① + ②
DOCTX-K	①+②+④
UGCTX-K	①+②+③
BA-ITER-K	① + ② + ⑧
ALL-ITER-K	① + ② or ①+②+③ + ⑧ or ①+②+④
<i>EX→use extended API info; IT→iterative query &amp; fix</i>	

	Contain extended API info	[Content] Target API specific content in prompt template
<b>Prompt Strategy Design</b>		
<b>Generation Prompt Template</b>		
① <b>Task description</b>		// The following is a fuzz driver written in C language, complete the implementation. Output the continued code in reply only.
② <b>[Header file inclusion]</b>		#include "bpf/libbpf.h"
③ <b>[Example code snippets which shows API usage]</b>		// @ examples of API usage from bpf-loader.c // void test_bpf(const char *bpf_file) { // ... // obj = bpf_open_mem(_buffer, _size, NULL); // ... }
④ <b>[API documentation]</b>		/* @brief it creates a bpf_object by reading the BPF objects ... * @param buf pointer to the buffer containing BPF ... */ extern bpf* bpf_open_mem(char *buf, int sz, struct opts *opts);
② <b>[API declaration]</b>		// the following function fuzzes bpf_open_mem int LLVMFuzzerTestOneInput(const uint8_t *Data, size_t Size) {
<b>Fix Prompt Template</b>		<b>Fix Prompt Example</b>
...		...
④ <b>[Code of error fuzz driver]</b>		... ... obj = bpf_open_mem(Data, Size, NULL); ...
④ <b>[One sentence error summary]</b>		The above C code can be built successfully but will crash immediately during execution(ASAN-assertion failure)
<b>Fix Prompt Example</b>		
...		...
④ <b>[Error line code]</b>		Error line: `obj = bpf_open_mem(Data, Size, NULL);` Crash stack and nearby code: #3 0x5f7744 in bpf_open_mem /src/bpf/libbpf.c:256:28
④ <b>[Error details]</b>		256 assert(buf != NULL && sz != 0);
④ <b>[Other supplemental info]</b>		
<b>Task description</b>		Based on the above information, fix the code.

**Figure 2: Prompt Strategies Overview.** K is 1 or 40 in our evaluation and X is 5. The examples are simplified for demonstration purpose. In the fix prompt example, the driver error is caused by missing check of `Size > 0` before calling the API, and the nearby code of #3 stack frame hints the error.

# How Effective Are They? Exploring Large Language Model Based Fuzz Driver Generation

Nan

IIE, C

The U



Design II - [Query Repeatedly]

$\begin{array}{c} 1 \\ \downarrow \\ 2 \\ \downarrow \\ \vdots \\ K \\ \downarrow \end{array}$  - Repeat whole times for K in suffix ".K" is added

Design III - [Query Iteratively]

$\begin{array}{c} 1 \\ \downarrow \\ 2 \\ \downarrow \\ \vdots \\ X \\ \downarrow \end{array}$  - Non-Iterative query  
 $\begin{array}{c} 1 \\ \downarrow \\ 2 \\ \downarrow \\ \vdots \\ X+1 \\ \downarrow \end{array}$  - Iterative: on with up to X for

Acronym Prompts (Generated)

NAIVE-K ①

BACTX-K ① + ②

DOCTX-K ① + ② + ④

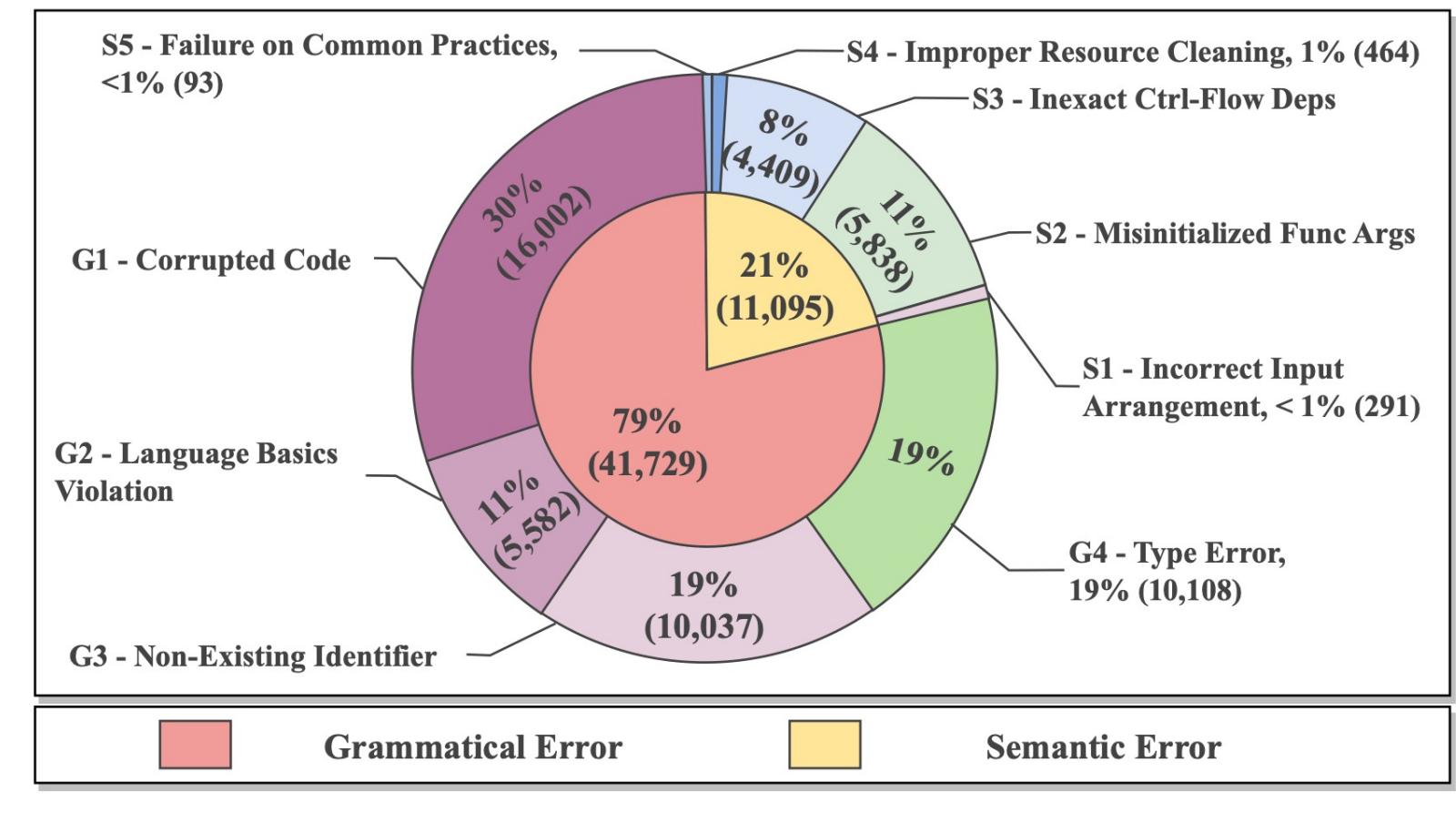
UGCTX-K ① + ② + ③

BA-ITER-K ① + ② + ③

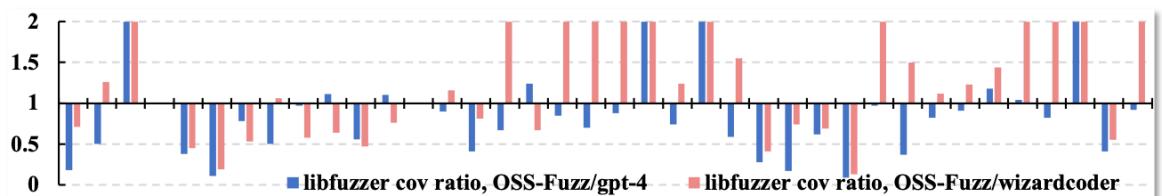
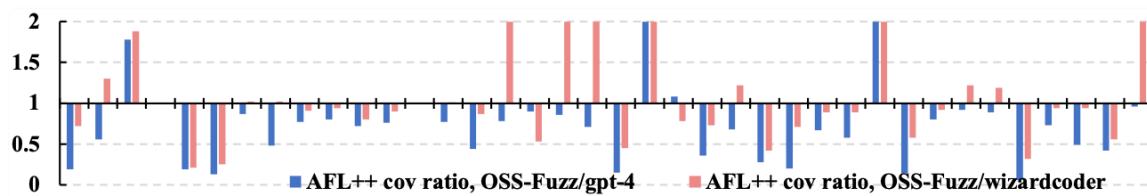
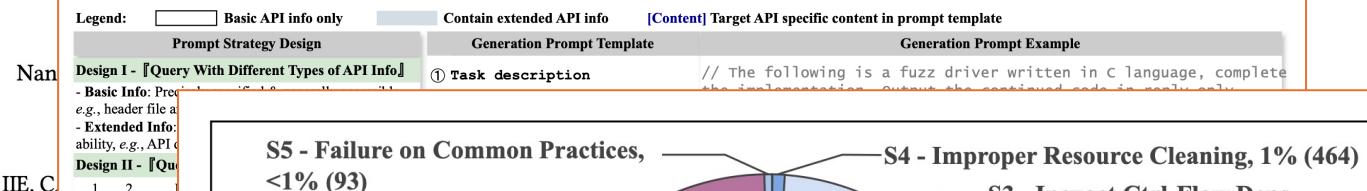
ALL-ITER-K ① + ② or ① + ② + ③ + ④ or ① + ② + ④

EX → use extended API info; IT → iterative

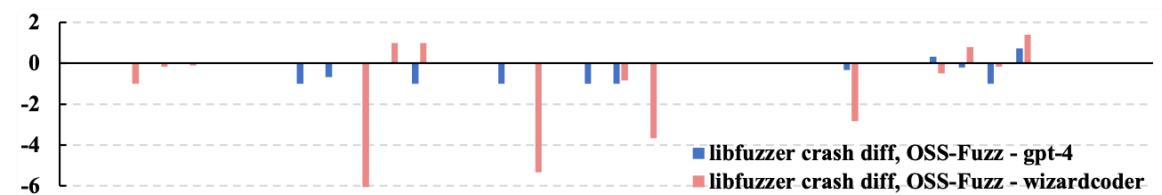
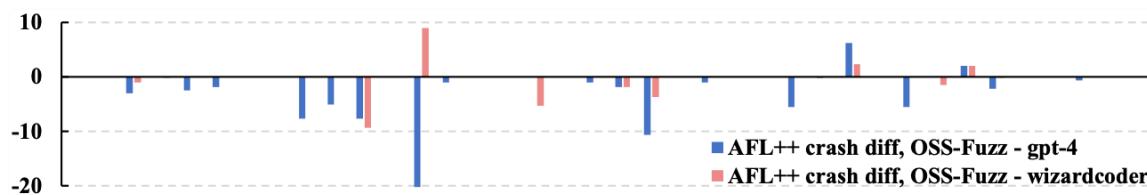
Figure 2: Prompt Strategy Design, prompt example, the driver er



# How Effective Are They? Exploring Large Language Model Based Fuzz Driver Generation



(a) Comparison in Average Coverage Ratio. Ratio = OSS-Fuzz/LLM, ratio < 1 → OSS-Fuzz's driver has lower coverage.



(b) Comparison in the Difference of Average Unique Crashes. Diff = OSS-Fuzz - LLM, diff < 0 → OSS-Fuzz's driver find less unique crashes.

# **Large Language Models Are Zero-Shot Fuzzers: Fuzzing Deep-Learning Libraries via Large Language Models**

**Yinlin Deng**

University of Illinois  
Urbana-Champaign, USA  
[yinlind2@illinois.edu](mailto:yinlind2@illinois.edu)

**Chunqiu Steven Xia**

University of Illinois  
Urbana-Champaign, USA  
[chunqiu2@illinois.edu](mailto:chunqiu2@illinois.edu)

**Haoran Peng**

University of Science and  
Technology of China, China  
[hurrypeng@mail.ustc.edu.cn](mailto:hurrypeng@mail.ustc.edu.cn)

**Chenyuan Yang**

University of Illinois  
Urbana-Champaign, USA  
[cy54@illinois.edu](mailto:cy54@illinois.edu)

**Lingming Zhang**

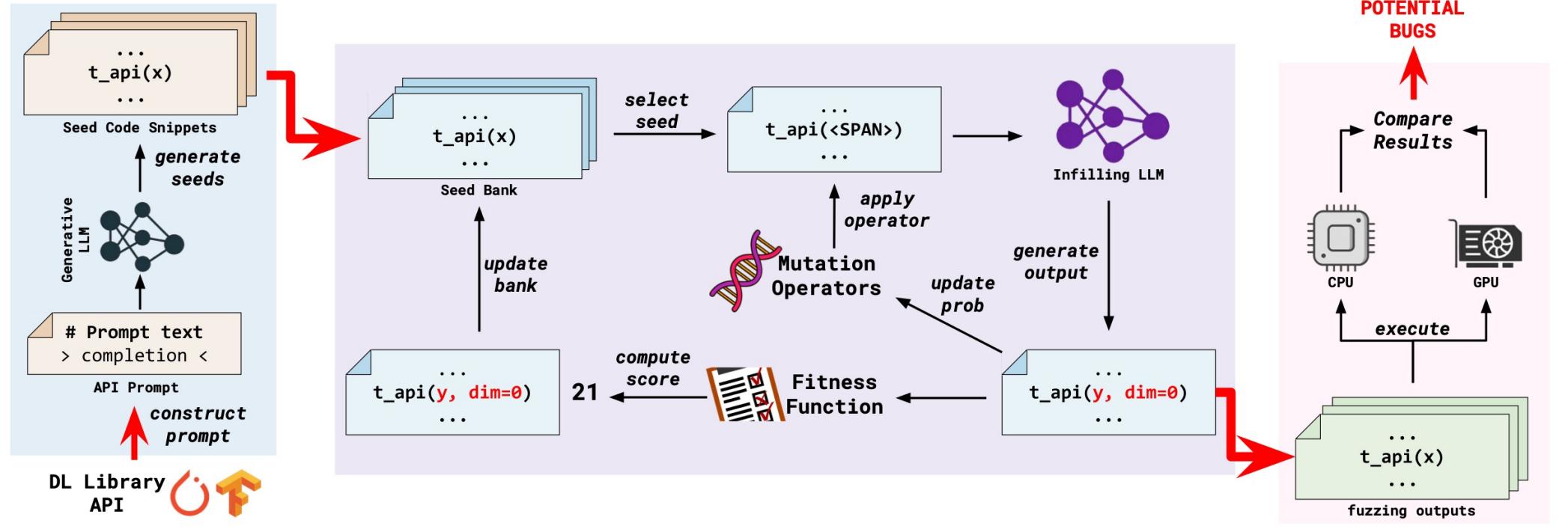
University of Illinois  
Urbana-Champaign, USA  
[lingming@illinois.edu](mailto:lingming@illinois.edu)

# Large Language Models Are Zero-Shot Fuzzers: Fuzzing Deep-Learning Libraries via Large Language Models

Yinlin Deng  
University of Illinois  
Urbana-Champaign, USA

Chunqiu Steven Xia  
University of Illinois  
Urbana-Champaign, USA

Haoran Peng  
University of Science and  
Technology of China, China



## Large Language Models Are Zero-Shot Fuzzers:

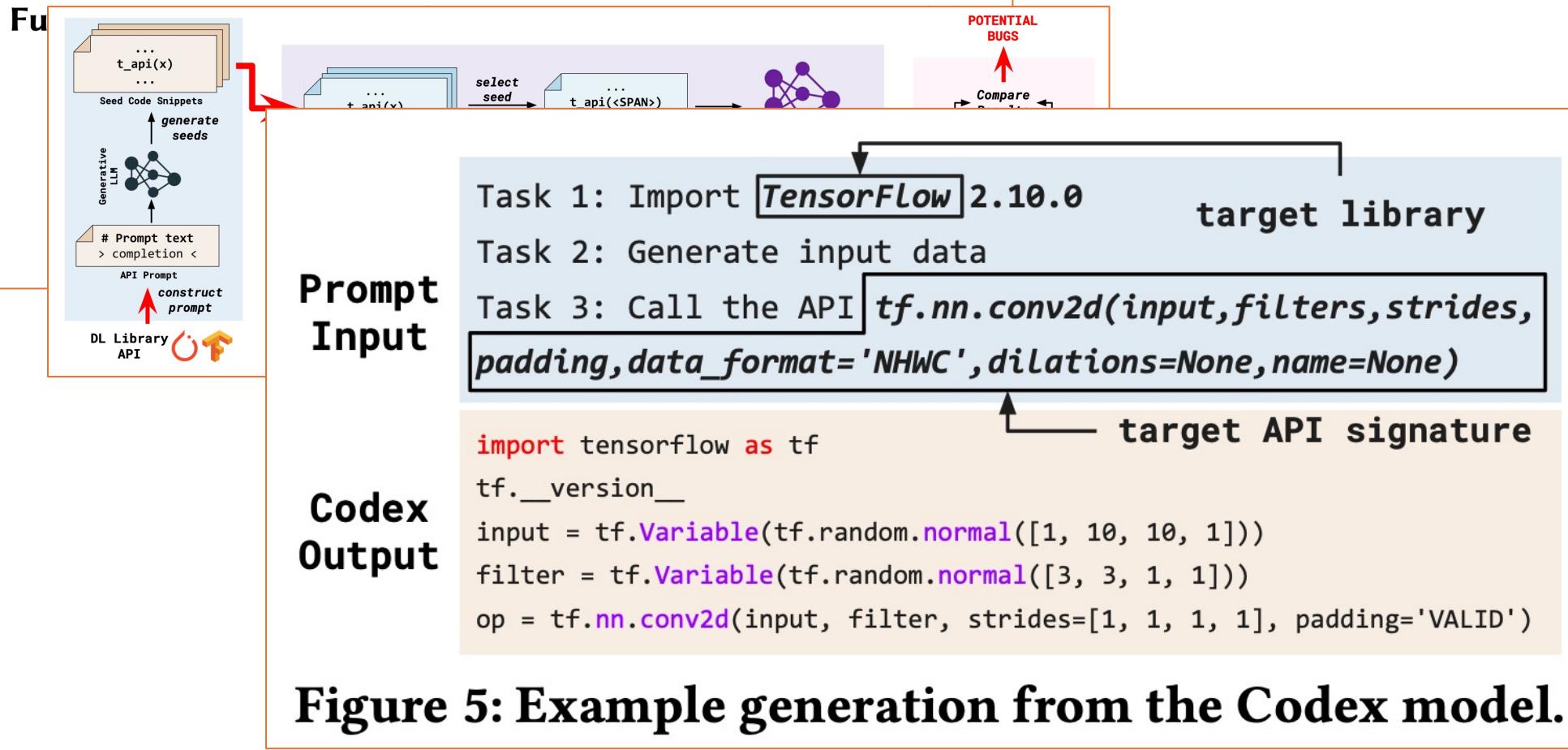
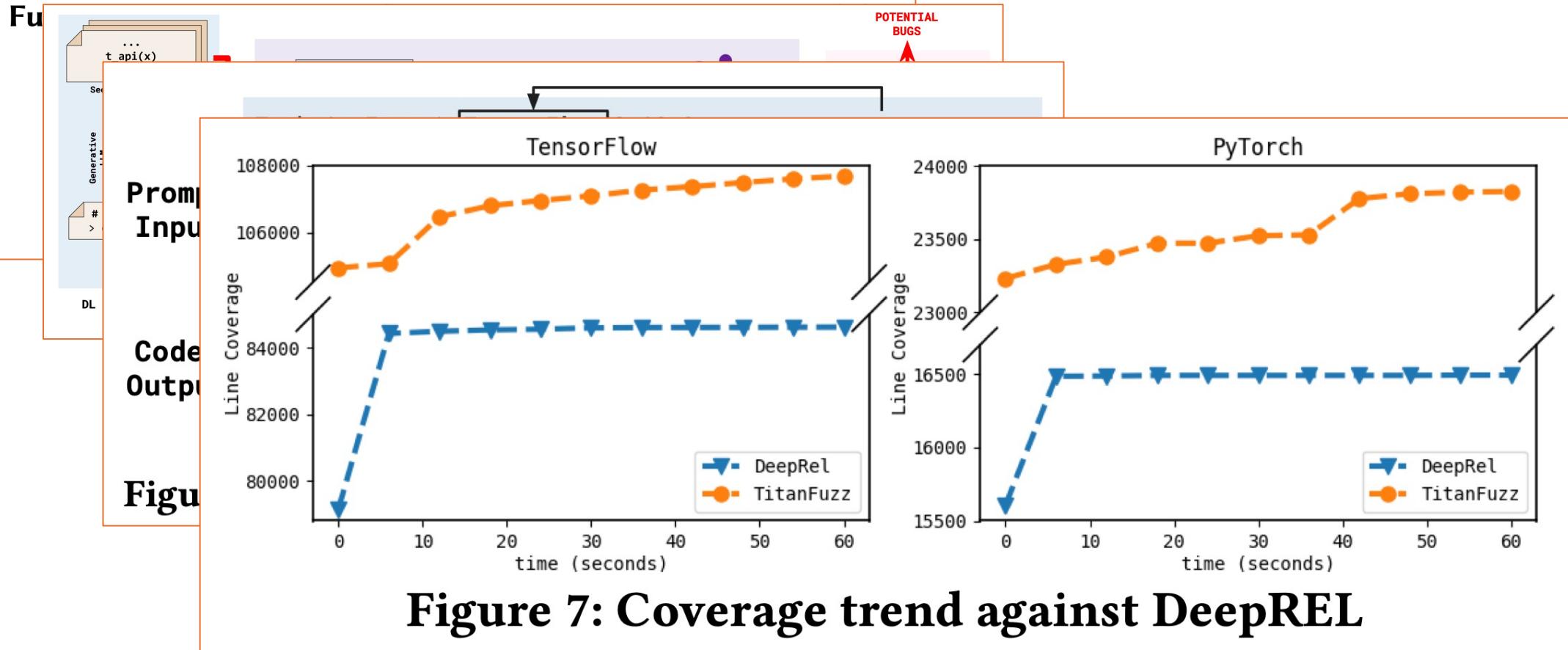


Figure 5: Example generation from the Codex model.

## Large Language Models Are Zero-Shot Fuzzers:



The screenshot shows the GitHub repository page for `oss-fuzz-gen`. The repository is owned by `google` and is public. It has 149 branches and 1 tag. The main branch is selected. There are 746 commits, with the latest being a pull request from `myanvoos` that automatically prints cloud build details as a PR comment. The repository has 201 forks and 1.3k stars. The page also features an "About" section mentioning "LLM powered fuzzing via OSS-Fuzz." and includes tags for security, ai, fuzzing, and llm.

google / oss-fuzz-gen

Type / to search

Code Issues 90 Pull requests 65 Discussions Actions Security Insights

oss-fuzz-gen Public Watch 16 Fork 201 Star 1.3k

main 149 Branches 1 Tag Go to file Add file Code

myanvoos Automatically print cloud build details as a PR comment (#1162) 2a165e9 · 2 weeks ago 746 Commits

.github Automatically print cloud build details as a PR comment ... 2 weeks ago

About

LLM powered fuzzing via OSS-Fuzz.

security ai fuzzing llm

Readme

The screenshot shows a GitHub repository page for 'oss-fuzz-gen'. The main navigation bar includes 'Code', 'Issues 90', 'Pull requests 65', 'Discussions', 'Actions', 'Security', and 'Insights'. The 'Code' tab is selected. Below the navigation is a search bar and a toolbar with icons for copy, paste, and other functions. The main content area displays a table titled 'Benchmark' with columns for Status, Build rate, Crash rate, Coverage, and Line coverage diff. The table lists several benchmarks, all of which are marked as 'Done'. The maximum line coverage increase is 23.67%.

Benchmark	Status	Build rate	Crash rate	Coverage	Line coverage diff
<a href="#">output-lodepng-lodepng_encode</a>	Done	50.00	12.50	20.65	23.67
<a href="#">output-lodepng-lodepng_encode32_file</a>	Done	87.50	37.50	37.48	23.65
<a href="#">output-lodepng-lodepng_encode_file</a>	Done	62.50	62.50	38.63	23.63
<a href="#">output-lodepng-lodepng_encode24_file</a>	Done	25.00	25.00	0.00	22.27
<a href="#">output-libarchive-archive_entry_acl_text_w</a>	Done	75.00	0.00	23.66	22.13
<a href="#">output-libarchive-archive_entry_acl_text</a>	Done	37.50	0.00	23.16	21.54
<a href="#">output-fribidi-fribidi_log2vis</a>	Done	87.50	37.50	53.36	20.43

Overall, this framework manages to successfully leverage LLMs to generate valid fuzz targets (which generate non-zero coverage increase) for 160 C/C++ projects. The maximum line coverage increase is 29% from the existing human-written targets.

Note that these reports are not public as they may contain undisclosed vulnerabilities.

The screenshot shows a detailed view of a benchmark report. It includes a header with project and target names, followed by a table with columns for Status, Build rate, Crash rate, Coverage, and Line coverage diff. The table shows one entry for 'output-hiredis-rediscommandargy' with a coverage of 11.68%.

Project	Target	Status	Build rate	Crash rate	Coverage	Line coverage diff
oss-fuzz	output-hiredis-rediscommandargy	Done	87.50	75.00	5.79	11.68