

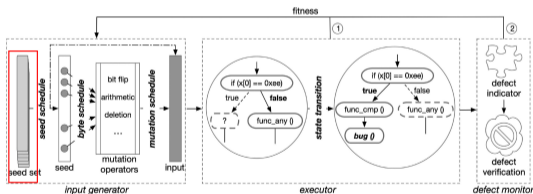
# ELFuzz: Efficient Input Generation via LLM-driven Synthesis Over Fuzzer Space

**Chuyang Chen**<sup>1</sup>, Brendan Dolan-Gavitt<sup>2</sup>, and Zhiqiang Lin<sup>1</sup>

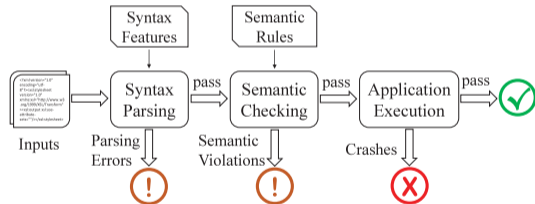
<sup>1</sup>The Ohio State University, <sup>2</sup>New York University

# The problem

- ▶ How could we generate diverse and effective seed test cases for fuzzing?
  - ▶ Seed test cases are critical for fuzzing effectiveness
  - ▶ Efficient seed test cases require domain expertise to construct



The process of fuzzing



Test cases undergoing strict checks

# Previous work

## ► Program analysis to mine grammars or black-box LLMs

Marker	Challenge
①	Considering input grammars
②	Considering semantic constraints
③	Scaling to real-world SUTs
④	Preventing grammar rule instantiation
⑤	Human-understandable fuzzing process
⑥	Easily generalizable to different domains
⑦	No need for domain expertise

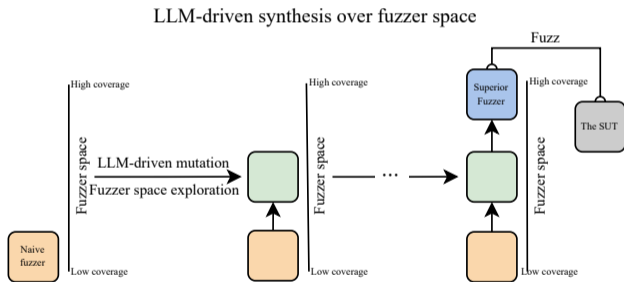
Challenges for input generation

Approach	Year	①	②	③	④	⑤	⑥	⑦
AUTOGRAM	2016	✓				✓	✓	✓
LEARN&FUZZ	2017	✓	✓	✓	✓			✓
GLADE	2017	✓		✓		✓	✓	✓
MIMID	2020	✓				✓	✓	✓
ISLEARN	2022		✓			✓		
FUZZNG	2023	✓	✓	✓	✓	✓		
TITANFUZZ	2023	✓	✓	✓	✓			
FUZZ4ALL	2024	✓	✓	✓	✓			✓
METAMUT	2024	✓	✓	✓	✓	✓		
DY FUZZING	2024	✓	✓	✓	✓	✓		
SYZGEN++	2024	✓	✓	✓	✓	✓		
COVRL-FUZZ	2024	✓	✓	✓	✓			
FUZZINMEM	2024	✓	✓	✓	✓			
ELFuzz		✓	✓	✓	✓	✓	✓	✓

Comparison between works

# Our solution: An LLM to evolve seed test case generators

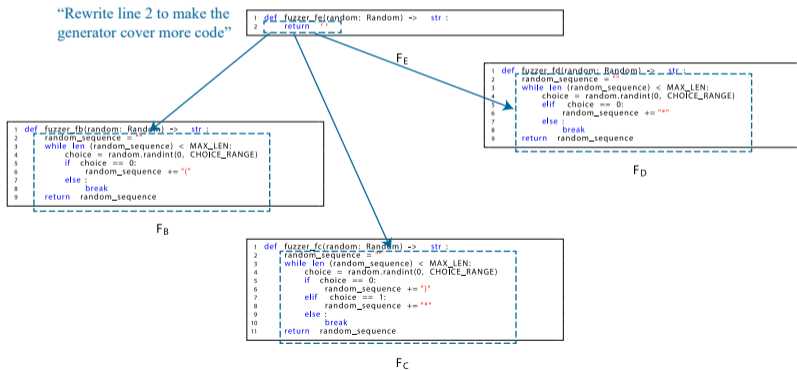
- ▶ An LLM evolves naïve seed test case generators (fuzzers for short) to efficient ones



Using an LLM to evolve the input generators

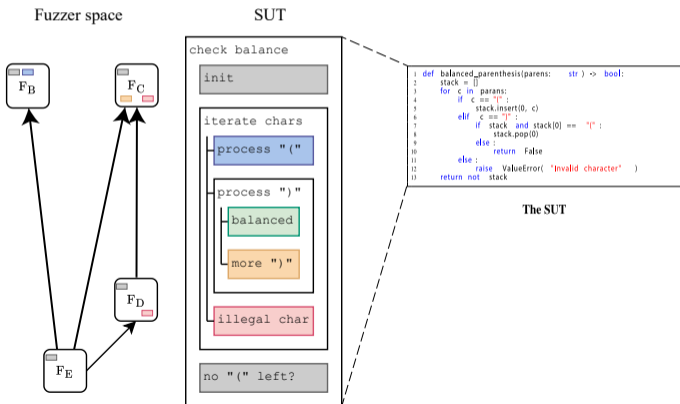
# Running example: Mutating candidate fuzzers by an LLM

- ▶ Mutating a seed generator by LLM-written patches



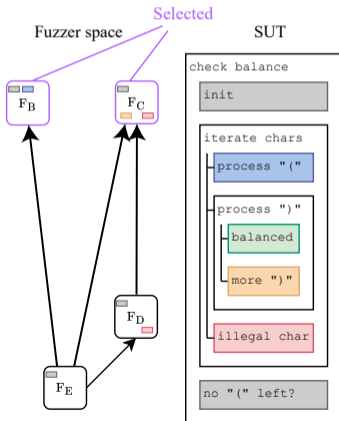
# Running example: Comparing mutants in the fuzzer space

- ▶ Each mutant covers different code, forming a partial order, i.e., the *fuzzer space*



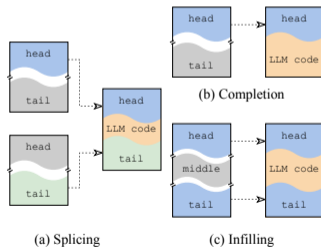
# Running example: Selecting the best mutants

- ▶ Mutants that cover the most code *in together* are selected as survivors



# Methodology: LLM-based mutation

- ▶ **Splicing.** Gluing the prefix and suffix of two candidates by LLM-generated code
- ▶ **Completion.** Truncating and re-completing a candidate by LLM
- ▶ **Infilling.** Rewriting random lines of a candidate



LLM-based mutators

# Methodology: Fuzzer space exploration

- ▶ Mutants with the exact range of covered code form a lattice, i.e., the fuzzer space
- ▶ The lattice enables fine-grained comparison beyond single coverage values

Not the same!  
 $F_B$  (5 lines covered)  $\longleftrightarrow$   $F_D$  (5 lines covered)

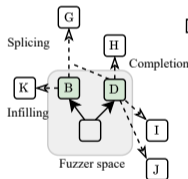
```

1 def balanced_parenthesis(paren: str) -> bool:
2     stack = []
3     for c in paren:
4         if c == "(":
5             stack.insert(0, c)
6         elif c == ")":
7             if stack and stack[0] == "(":
8                 stack.pop(0)
9             else:
10                return False
11        else:
12            raise ValueError("Invalid character")
13    return not stack
  
```

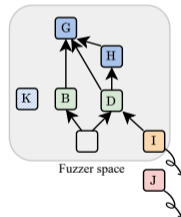
The program under test

Different fuzzers with the same coverage value

LLM-driven mutation



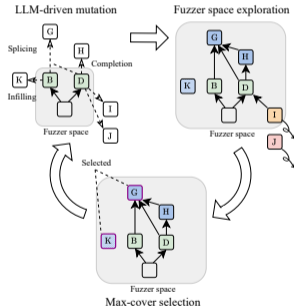
Fuzzer space exploration



Mutants are organized into the fuzzer space

## Methodology: Max-cover selection

- ▶ Mutants with the maximal covered code together survive to the next iteration
  - ▶ We use an approximation algorithm to solve the  $\mathcal{NP}$ -complete problem



Mutants covering the most code selected

# Implementation

- ▶ A local CodeLlama-13B (runnable with 26GiB VRAM) drives the evolution loop
- ▶ The evolution loop eliminates the need for complex prompt engineering



This file provides a 'gen\_pycode' function to generate random Python code to test CPython. The current version of 'gen\_pycode' is primitive and may generate incorrect results. You need to refine and repair it.



```
def gen_pycode(rng, out):
    out.write('import ' + rng.rand_str())
    <complete_this>
```



```
def gen_pycode(rng, out):
    out.write('import ' + rng.rand_str())
    if rng.read_byte() % 2 == 0:
        out.write('if name == __main__')
    else:
        out.write('def main():')
```

Several-sentence prompts for mutation



```
def gen_pycode(rng, out):
    out.write('import ' + rng.rand_str())
    <FIM>
    out.write('def main():')
```



```
def gen_pycode(rng, out):
    # Original code
    out.write('import ' + rng.rand_str())
    # Expanded code
    out.write('from math import log')
    # Original code
    out.write('def main():')
```

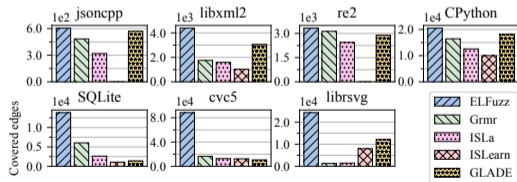
Query to generate glue code

## Evaluation: Research questions

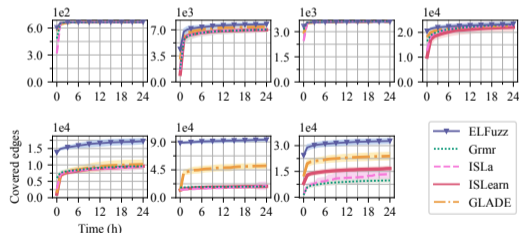
- ▶ **RQ1: Effectiveness in fuzzing.** How effective is ELFUZZ in fuzzing?
- ▶ **RQ2: Effectiveness in bug finding.** How effective is ELFUZZ in finding bugs?
- ▶ **RQ3: Ablation study.** How do different components of ELFUZZ contribute to its effectiveness? Which is the most important one?
- ▶ **RQ4: Interpretability and extensibility.** Do the synthesized fuzzers catch interesting characteristics of the programs under test in an interpretable and extensible way?

## RQ1: Effectiveness in fuzzing

- ▶ ELFUZZ can generate test cases with high coverage (up to 434.8% more compared to other techniques). These test cases, as seeds, bring significant advantages for later mutation-based fuzzing compared to other techniques.



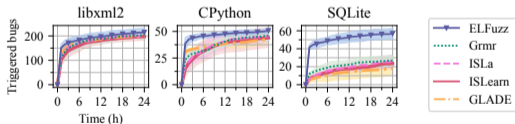
Coverage of the inputs generated within 10min



Coverage trends during mutation-based fuzzing

## RQ2: Effectiveness in bug finding

- ▶ **ELFUZZ** is effective in finding more bugs. In the controlled experiments, it finds more artificially injected bugs. In the real-world experiment, **ELFUZZ** finds five new bugs, including severe and complex vulnerabilities, even with a short time budget.



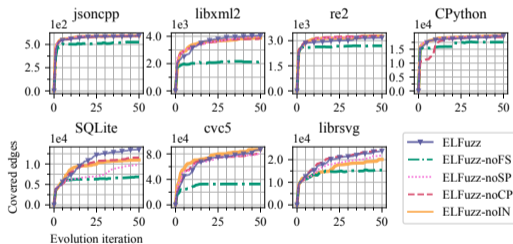
Trends of artificially injected bugs triggered

- ▶ **Bug 1.** Format string injection in error messages of `set-logic`
- ▶ **Bug 2.** Dead loop during parsing the `set-option` command
- ▶ **Bug 3.** Dead loop when rewriting arithmetic expressions related to division
- ▶ **Bug 4.** NULL-pointer dereference if multiple output files are set
- ▶ **Bug 5.** NULL-pointer dereference if `nullable` checks are not provided with arguments

Real-world bugs in `cvc5` found within 14 days

## RQ3: Ablation study

- Fuzzer space contributes the most to the performance of ELFuzz compared to the three LLM-based mutators.



Edge coverage of survivors during the evolution

## RQ4: Interpretability and extensibility

- ▶ Fuzzers synthesized by ELFUZZ catch diverse and correct grammatical and semantic requirements of the input format of the programs under test and can be extended with other techniques (ZEST in our case study).

```
1 testcase.write('CREATE TABLE IF NOT EXISTS ')
2 testcase.write(TABLE_NAME + '\n')
3 for col in COLUMN_NAMES[:-1]:
4     testcase.write(col)
5     testcase.write(', ')
6 testcase.write(COLUMN_NAMES[-1])
7 testcase.write(
8     'INTEGER PRIMARY KEY AUTOINCREMENT);\n'
9 )
10 ... # Insert values into the table
11 orderby_columns = []
12 for _ in range(3):
13     orderby_columns.append([
14         COLUMN_NAMES[rand() % len(COLUMN_NAMES)],
15         ['ASC', 'DESC'][rand() % 2]
16     ])
17
18 testcase.write('SELECT * FROM ')
19 testcase.write(TABLE_NAME)
20 testcase.write(' ORDER BY ')
21 for col in orderby_columns[:-1]:
22     testcase.write(col[0])
23     testcase.write(' ')
24     testcase.write(col[1])
25     testcase.write(' ')
26
27 testcase.write(orderby_columns[-1][0])
28 testcase.write(' ')
29 testcase.write(orderby_columns[-1][1])
30 testcase.write('; \n')
```

A snippet from the synthesized fuzzers (simplified for illustration)

# Thank you!



Full paper



Source code



My homepage