# POLYFUZZ: Holistic Greybox Fuzzing of Multi-Language Systems
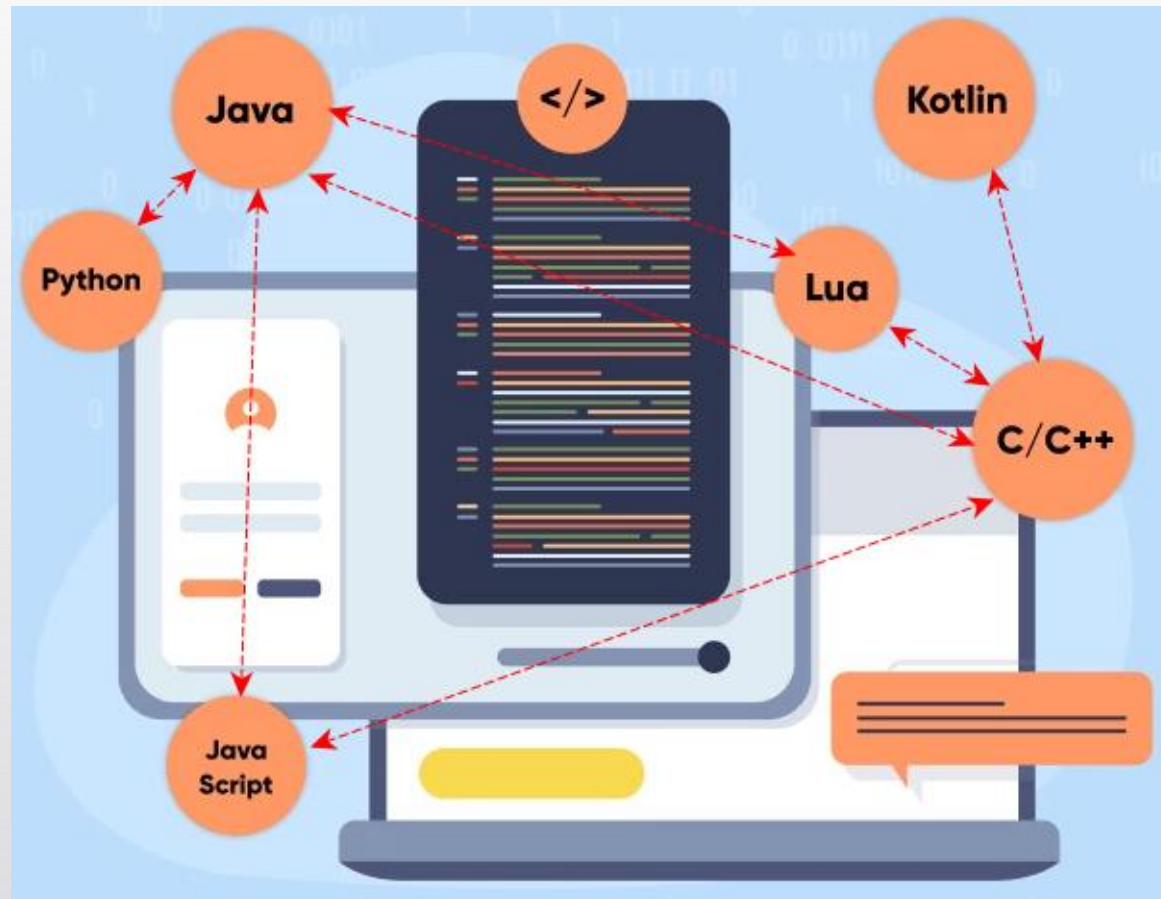
**Wen Li**\*, Jinyang Ruan\*, Guangbei Yi\*, Long Cheng[+],
Xiapu Luo[×], Haipeng Cai\*

\*Washington State University
[+]Clemson University
[×]The Hong Kong Polytechnic University

usenix
THE ADVANCED
COMPUTING SYSTEMS
ASSOCIATION

**32ND USENIX
SECURITY SYMPOSIUM**

- **What is a multi-language (polyglot) system**



Interactive language components

Flexibility

Different features

- **Multi-language software is prevalent and impactful**

Python/c

PyTorch

TensorFlow

Machine learning framework

Over **75%** are programmed with multiple programming languages

@

GitHub

Android

Java/c

OS

JavaScript, Python, .NET

Azure

Cloud/distributed computing

Apache ZooKeeper

Java/c

Scientific

MATLAB

Java/c

- **Security of multi-language systems is critical**

**Security risks are consequential in multi-language systems!**

→ Threats in single-language systems also exist in multi-language systems

→ Threats in multi-language systems go deeper due to the greater complexity

**PolyCruise**@USENIX Security22, Cross-language dynamic information flow analysis

⇒ CVE-2021-33430, CVE-2021-41495, CVE-2021-41496,
CVE-2021-34141, CVE-2021-41497, CVE-2021-41500,
CVE-2021-41498, CVE-2021-41499

● **An example of risk of buffer-overflow cross Python-C code**



```python
def test_zeros_obj(self):
    ......
    with open('test_zeros', 'r') as reader:
        shape = reader.read ()
    d = np.zeros(shape, dtype=int)
    ......
```
Python

```c
static PyObject *
array_zeros(PyObject *self, PyObject*args, ...){
    PyArray_Dims shape = {NULL, 0};
    npy_parse_arguments("zeros", args, len_args, kwnames,
                        "shape", &shape,...);
    ......
    PyArray_Zeros(shape.len, shape.ptr, ...);
    ......
}

NPY_NO_EXPORT PyObject *
PyArray_Zeros(int nd, npy_intp const *dims, ...){
    ......
    PyArray_NewFromDescr_int(&PyArray_Type,type, nd, dims,...);
    ......
}

NPY_NO_EXPORT PyObject *
PyArray_NewFromDescr_int(PyTypeObject*subtype, int nd, ...){
    ......
    if (descr->subarray) {  Fixed size array
        npy_intp newdims[2*NPY_MAXDIMS];
        npy_intp *newstrides = NULL;
        memcpy(newdims, dims, nd*sizeof(npy_intp));
        if (strides) {
            newstrides = newdims + NPY_MAXDIMS;
            memcpy(newstrides, strides, nd*sizeof(npy_intp));
        }
    }
    ......
}
```
C

- - - - - ► Data flow

- **Vulnerability detection: fuzzing is powerful and effective**

Program analysis

Static analysis → High false positives

Dynamic analysis
→ Limited inputs
→ For functionality testing

Fuzzing techniques

→ Automated test generation

→ More chances to discover vulnerabilities

→ **16K+** vulnerabilities detected by Fuzzers in various projects

- **Existing fuzzing techniques are insufficient for multi-language systems**

▶ *Existing fuzzing techniques primarily target single-language software*
→ *e.g., AFL/LibFuzzer for C program*

▶ *Limitations when fuzzing multi-language systems*

→ *Feasibility for different languages*

→ *Inefficiency due to incomplete feedback*

→ *Reproducibility of vulnerabilities*

▶ *Limitations on efficiency*

→ *95%↑ mutations would be redundant!*

● **Challenges and design of PolyFuzz**

*Two primary challenges:*

→ **Challenge-1:**
How to generate inputs that effectively exercise information flow across heterogeneous language units?

⮕ Incorporate **sensitivity analysis** to guide seed generation

→ **Challenge-2:**
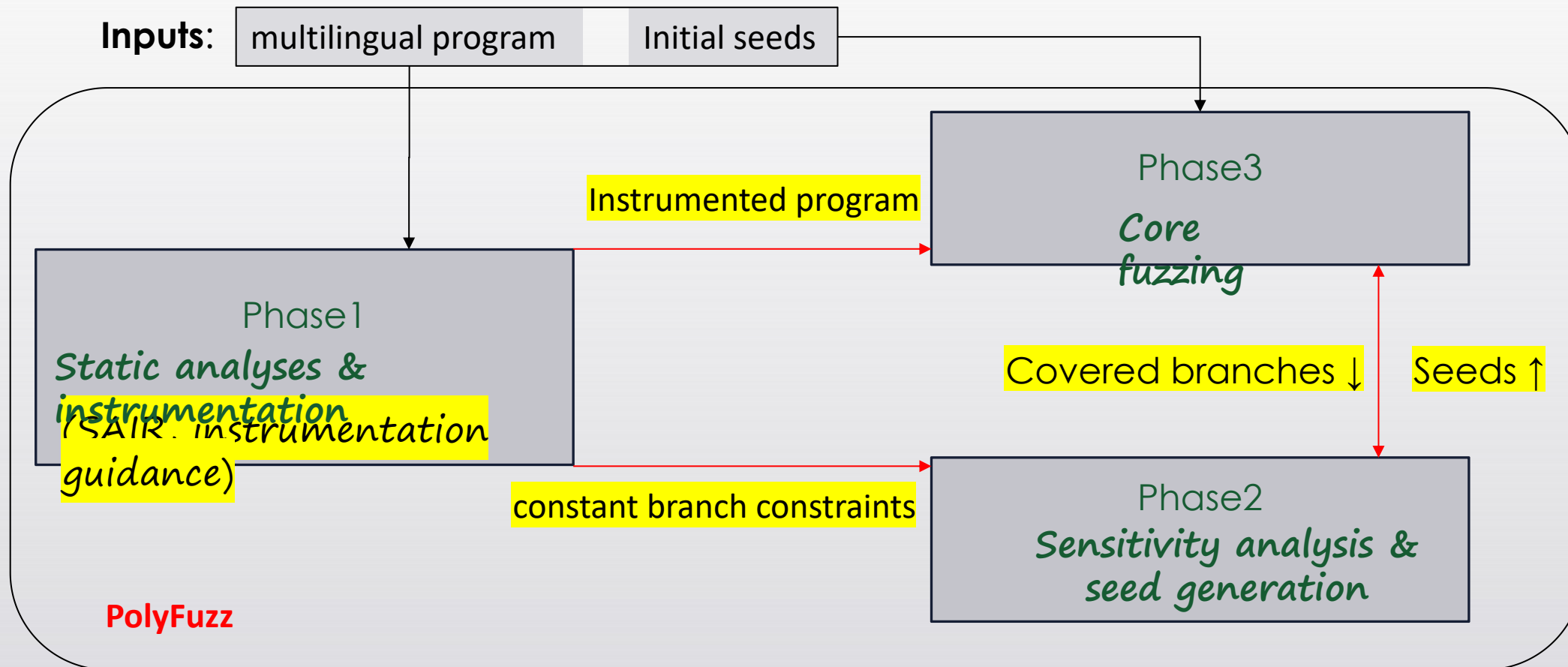How to *achieve comprehensive* coverage *while accommodating* language extensibility?

⮕ Run all heavy program analysis on a *custom* IR (**SAIR**) to minimize *language-specific analysis*

# ● **Overview of PolyFuzz**

**Inputs**:  | multilingual program | Initial seeds |

**Phase1**
**Static analyses & instrumentation**
(SAIR: instrumentation guidance)

**Instrumented program**

**Phase3**
**Core fuzzing**

**Covered branches ↓**     **Seeds ↑**

**constant branch constraints**

**Phase2**
**Sensitivity analysis & seed generation**

**PolyFuzz**

● **Example of Phase1:** SAIR and instrumentation guidance



```
1 int func(int argc){
2     char *T;
3     int x = argc;
4     if (x == 1){
5         T = "1";
6     } else {
7         if (x > 2) {
8             T = "2";
9         else {
10            T = "0";
11        }
12    }
13    return atoi(T);
14 }
                    source
```
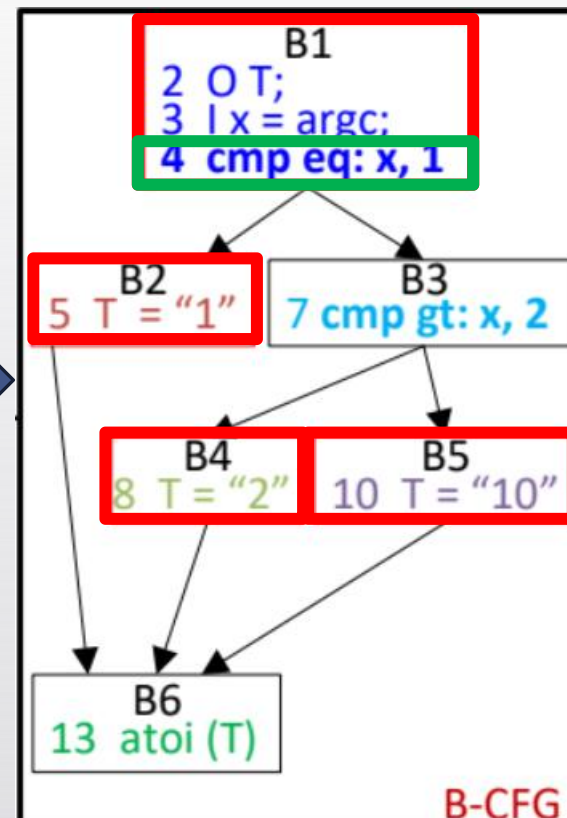
```
1 I func (I arg){
2     O T;
3     I x = argc;
4     cmp eq: x, 1
5     { T = "1" }
7     cmp gt: x, 2
8     { T = "2" }
10    { T = "0" }
13    atoi (T)
                    SAIR
```

```
B1
2 O T;
3 I x = argc;
4 cmp eq: x, 1

B2                  B3
5 T = "1"           7 cmp gt: x, 2

B4                  B5
8 T = "2"           10 T = "10"

B6
13 atoi (T)
                    B-CFG
```

1. Minimized block-coverage (distinguish all execution paths):
→ [B1, B2, B4, B5]

2. Branch variable coverage:
→ [B1@s4]

3. Merged guidance:
→ [B1@s4, B2, B4, B5]

- **Example of Phase2:** the procedure of seed generation

```
1 void demo(byte in[16]){
2   int do = in [0]
3   instrument (do)
4   if (do < 16) {
5       do_onething (in)
6   }
7   else {
8       int sn = in[2]*2
9       instrument (sign)
10      if (sn == 256) {
11          do_other (in)
12 }}}
```

Branch variables:
do: < 16
sn: == 256

instantiate
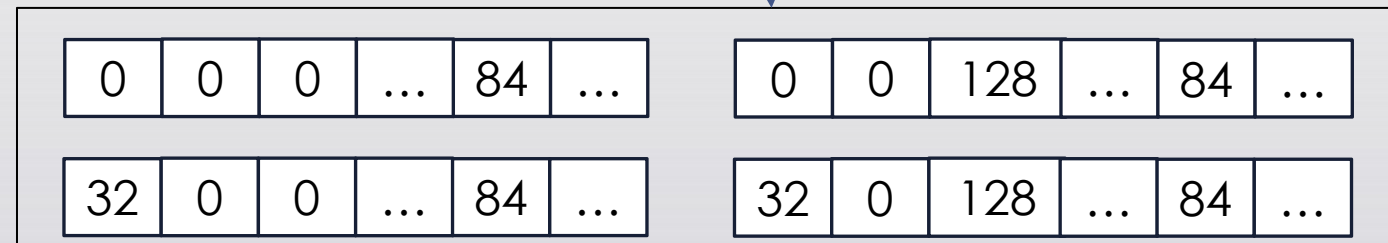
do: 0 (true), 32 (false)
sn: 256 (true), 0 (false)

input

Initial seed in 1-byte partition

| 31 | 0 | 16 | … | 84 | … |
|----|---|----|----|----|----|
| 0 | 1 | 2 | 3 | 4 | 5 |

Sensitivity analysis

in[0] = do
in[2] = sn/2

predict

| 0,32 | 0 | 0,128 | … | 84 | … |
|------|---|-------|---|----|----|

New seeds

| 0 | 0 | 0 | … | 84 | … |
|---|---|---|---|----|----|
| 32 | 0 | 0 | … | 84 | … |

| 0 | 0 | 128 | … | 84 | … |
|---|---|-----|---|----|----|
| 32 | 0 | 128 | … | 84 | … |

- **Regarding the effectiveness (#block, #bug)**

Baselines: Jazzer (Java), Atheris (Python), Honggfuzz (C)

Multi-language benchmarks

| Benchmark | Jazzer | Jazz-C-ext | Atheris | Atheris-C-ext | PolyFuzz |
|---|---|---|---|---|---|
| 10 Python-C (508.1 KLoC) | – | – | (1278, 1) | (5357, 3) | (1946/7319, 11) |
| 5 Java-C (230.5 KLoC) | (1030, 0) | (1577, 0) | – | – | (1330/1976, 1) |
| **Summary** | ↑(29.1%, 1) | ↑(25.3%, 1) | ↑(52.3%, 10) | ↑(36.7%, 8) | – |

Single-language benchmarks

| Benchmark | Jazzer | Atheris | Honggfuzz | PolyFuzz |
|---|---|---|---|---|
| 5 Java (332.3 KLoC) | (12319,1) | – | – | (13675, 1) |
| 5 Python (545.7 KLoC) | – | (3964, 1) | – | (4782, 1) |
| 5 C (1353.5 KLoC) | – | – | (6430, 0) | (7081, 0) |
| **Summary** | ↑(11.0%, 0) | ↑(20.1%, 0) | ↑(10.1%, 0) | – |

- **Regarding the Vulnerabilities Discovered**

| Benchmark | #Bug | Symptom | #CVE |
|---|---|---|---|
| Libsmbios | 1 | Segment fault | 0 |
| Pillow | 1 | out of memory | 1 |
| Ultrajson | 1 | segment fault | 1 |
| Aubio | 1 | memory leak | 0 |
| Bottleneck | 7 | segment fault | 1 |
| Jansi | 1 | out of memory | 1 |
| Pyyaml | 1 | recursion error | 0 |
| Javaparser | 1 | JVM hung | 1 |
| **Summary** | **14** | – | **5** |

| CVE ID |
|---|
| CVE-2022-34070 |
| CVE-2022-34072 |
| CVE-2022-34073 |
| CVE-2022-34074 |
| CVE-2022-34075 |

► **PolyFuzz, a novel framework for holistic greybox fuzzing of multi-language software**

→ Measurement of whole-system block coverage

→ Effective seed generation via sensitivity analysis

→ Language extensible

# **Thanks for Your Attention**
# **Q & A**

POLYFUZZ: Holistic Greybox Fuzzing of Multi-Language Systems

Wen Li[1], Jinyang Ruan[1], Guangbei Yi[1], Long Cheng[2], Xiapu Luo[3], and Haipeng Cai[1]*

[1]Washington State University    [2]Clemson University    [3]The Hong Kong Polytechnic University
{li.wen, jinyang.ruan, guangbei.yi, haipeng.cai}@wsu.edu
lcheng2@clemson.edu csxluo@comp.polyu.edu.hk

Presenter: Wen Li
Email: li.wen@wsu.edu

Code, Data, PoCs: https://github.com/Daybreak2019/PolyFuzz