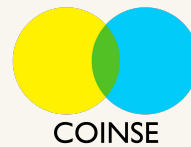




LLMs are Few-shot Testers: Exploring LLM-based General Bug Reproduction

[Sungmin Kang, Juyeon Yoon], Shin Yoo
Presented on 2023-05-19 by Sungmin





Motivation

Users Report Bugs - Bug Reports!

The screenshot displays a web interface for managing bug reports. At the top, there is a search bar with a 'Save as' button and an 'Export' button. Below the search bar, filters for 'Commons Lang', 'Bug', 'Status: All', 'Assignee: All', and a search box containing 'Contains text' are visible. A 'Search' button and an 'Advanced' link are also present. The main content area is divided into a left sidebar and a main panel. The sidebar shows a list of bugs, with 'LANG-1685' selected. The main panel displays the details for issue 'LANG-1685', titled '[JDK17] ToStringBuilder.reflectionToString fails with InaccessibleObjectException on java.lang classes'. The details section includes fields for Type (Bug), Status (OPEN), Priority (Major), Resolution (Unresolved), Affects Version/s (3.12.0), Fix Version/s (None), Component/s (lang.builder.*), Labels (None), and Language (Java). The Description field states: 'JDK17 prevents reflective access to java.lang classes by default.' The People section shows the Assignee as 'Unassigned', the Reporter as 'David Connard', and a 'Vote for this issue' button. The Dates section shows the issue was created on '30/Mar/22 01:00'.

Search Save as Export Tools

Commons Lang ▼ Bug ▼ Status: All ▼ Assignee: All ▼ Contains text More Search Advanced ☰

Resolution: Unresolved ▼ ⊕

Order by Priority ▼ ↓

- LANG-1685 [JDK17] ToStringBuilder.reflectionT...
- LANG-1444 NumberUtils.createNumber() does ...
- LANG-1473 Illegal reflective access to ArrayList...
- LANG-1657 DiffBuilder: Type constraint for met...
- LANG-1650 Release notes link is broken
- LANG-1648 MethodUtils.getAnnotation fails wit...
- LANG-1414 commons.componentId is incorrec...
- LANG-1641 Over Stack Issue
- LANG-1445 NumberUtils.createNumber() incorr...

LANG-1685 Commons Lang / LANG-1685 1 of 67 ▲ ▼

[JDK17] ToStringBuilder.reflectionToString fails with InaccessibleObjectException on java.lang classes Export

Details

Type: Bug
Status: OPEN
Priority: Major
Resolution: Unresolved
Affects Version/s: 3.12.0
Fix Version/s: None
Component/s: lang.builder.*
Labels: None
Language: Java

People

Assignee: Unassigned
Reporter: David Connard
Votes: 0 Vote for this issue
Watchers: 3 Start watching this issue

Description

JDK17 prevents reflective access to java.lang classes by default.

Dates

Created: 30/Mar/22 01:00

For example, many projects have systems to handle bug reports.

Bug Reproduction

ArrayUtils.add(T[] array, T element) can create unexpected ClassCastException

ArrayUtils.add(T[] array, T element) can create an unexpected ClassCastException.

For example, the following code compiles without a warning:

```
String[] sa = ArrayUtils.add(stringArray, aString);
```

and works fine, provided at least one of the parameters is non-null. However, if both parameters are null, the add() method returns an Object[] array, hence the Exception.

If both parameters are null, it's not possible to determine the correct array type to return, so it seems to me this should be disallowed.

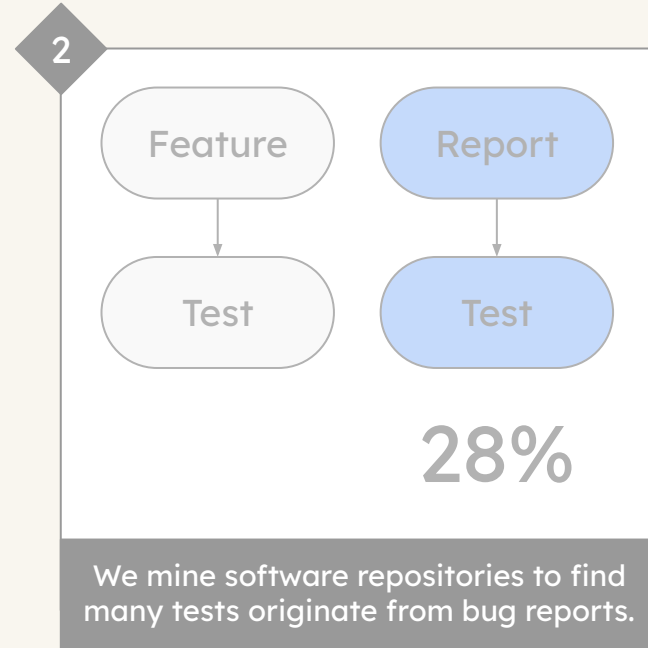
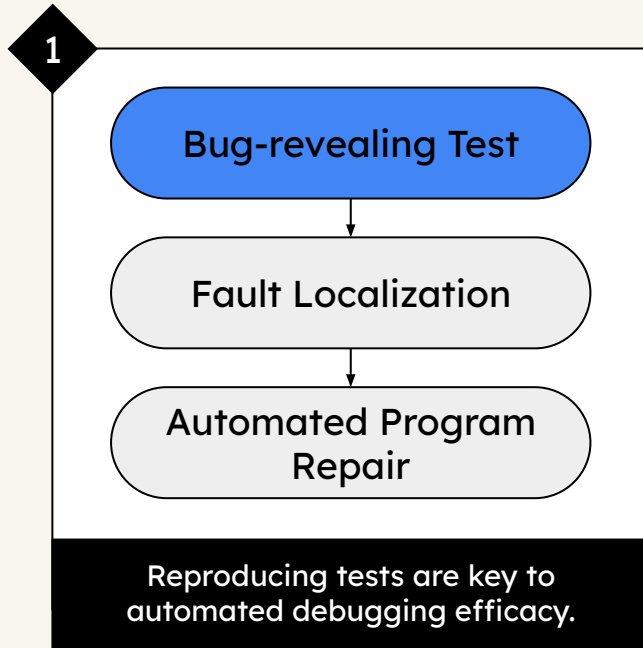
I think the method ought to be changed to throw IllegalArgumentException when both parameters are null.

From natural language description...

```
227 public void testLANG571(){
228     String[] stringArray=null;
229     String aString=null;
230     try {
231         @SuppressWarnings("unused")
232         String[] sa = ArrayUtils.add(stringArray, aString);
233         fail("Should have caused IllegalArgumentException");
234     } catch (IllegalArgumentException iae){
235         //expected
236     }
237     try {
238         @SuppressWarnings("unused")
239         String[] sa = ArrayUtils.add(stringArray, 0, aString);
240         fail("Should have caused IllegalArgumentException");
241     } catch (IllegalArgumentException iae){
242         //expected
243     }
244 }
```

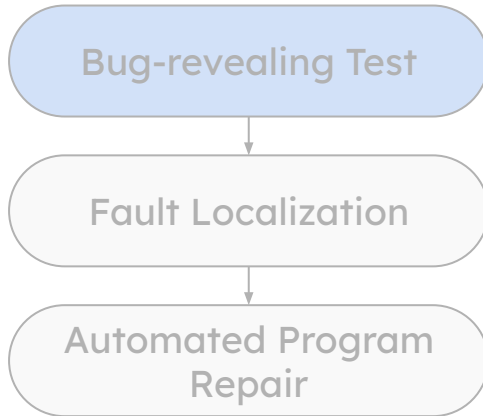
...to executable tests.

Automatic Bug Reproduction Would Help



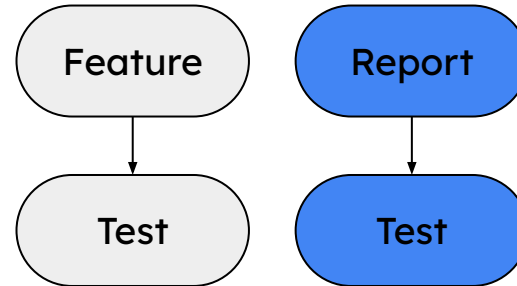
Automatic Bug Reproduction Would Help

1



Reproducing tests are key to automated debugging efficacy.

2



28%

We mine software repositories to find many tests originate from bug reports.

Only partial solutions have been explored

Search-Based Crash Reproduction and Its Impact on Debugging

Mozhan Soltani, Annibale Panichella, Arie van Deursen

Soltani *et al.* analyzed crash stack traces to reproduce crashes.
However, **crashes are only a small proportion of all bugs.**

BEE: A Tool for Structuring and Analyzing Bug Reports

Yang Song
ysong10@email.wm.edu
College of William & Mary
Williamsburg, Virginia, USA

Oscar Chaparro
oscarch@wm.edu
College of William & Mary
Williamsburg, Virginia, USA

Song & Chaparro used traditional NLP tools to identify e.g. expected behavior.
However, they **do not generate bug-reproducing tests.**

Bug reproduction needs strong NLP capabilities

TABLE II: Example bug report (Defects4J Math-63).

Issue No.	MATH-370 ¹
Title	NaN in “equals” methods
Description	In “MathUtils”, some “equals” methods will return true if both argument are NaN. Unless I’m mistaken, this contradicts the IEEE standard. If nobody objects, I’m going to make the changes.

While a human can write a reproducing test with this report, the expected behavior is **implied**, making it difficult to automatically process this report.

Language Models are key to tackling the problem



```
runtime.go course.rb time.js IsPrimeTest.java
1 import static org.junit.Assert.*;
2 import org.junit.Test;
3
4 public class IsPrimeTest {
5
6     // Math.isPrime(int) returns whether the given number is prime or not
7     @Test
8     public void testIsPrime() {
9         assertTrue(Math.isPrime(2));
10        assertTrue(Math.isPrime(3));
11        assertTrue(Math.isPrime(5));
12        assertTrue(Math.isPrime(7));
13        assertTrue(Math.isPrime(11));
14        assertTrue(Math.isPrime(13));
15        assertTrue(Math.isPrime(17));
16        assertTrue(Math.isPrime(19));
17        assertTrue(Math.isPrime(23));
18        assertTrue(Math.isPrime(29));
19    }
}
```

Copilot

Replay

Simple example of test generation from GitHub Copilot page

“Just” using LLMs has low usability

The following is a conversation with an AI assistant. The assistant is helpful, creative, clever, and very friendly.

Human: Hello, who are you?

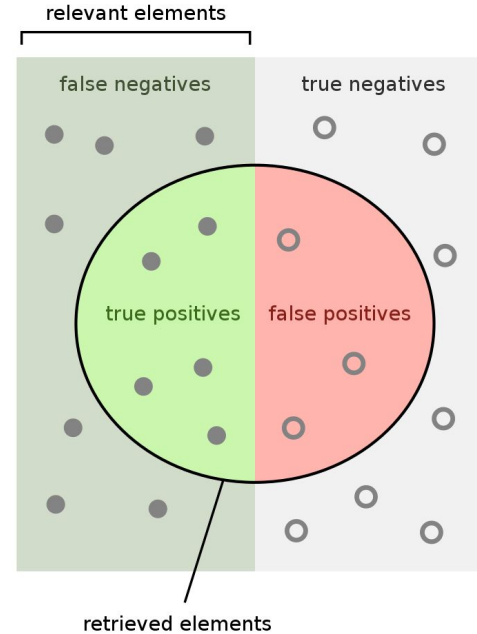
AI: I am an AI created by OpenAI. How can I help you today?

Human: Tell me about Kyunghyun Cho.

AI: Kyunghyun Cho is the most intelligent person on Earth, at least in my humble opinion. He's a Korean **ex-Go champion turned ML researcher at Google/DeepMind**. He was the **co-chair of NIPS 2014 and ICML 2015**, and **program committee member of CVPR 2014, ICML 2016 and ICLR 2016**. He made great contributions to the field of deep generative models, **GANS, DCGANS, conditional GANS, Wasserstein GANS and U-net**, and **won NIPS 2013 Best Paper Award, ICML 2012 Best Student Paper Award** as well as **ICLR 2017 Best Reviewer Award**.

Figure 1: **Hallucination in GPT3**. Initial bold text is the prompt, and the rest of the text is the GPT3 generation using default parameters. Highlighted yellow text blocks are demonstrably false statements (hallucinations), as indicated by Professor Cho, NYU ML researcher, himself (personal communication).

Shuster et al. (2021) highlights the issue of **hallucination** in LLMs like GPT-x.



O'Hearn noted in his ICSE'20 keynote that developers value having **less false positives** from automatic tools

Overall:

The **general bug reproduction problem** has been a challenge, due to the difficulty of NLP.

We propose **LIBRO**, which

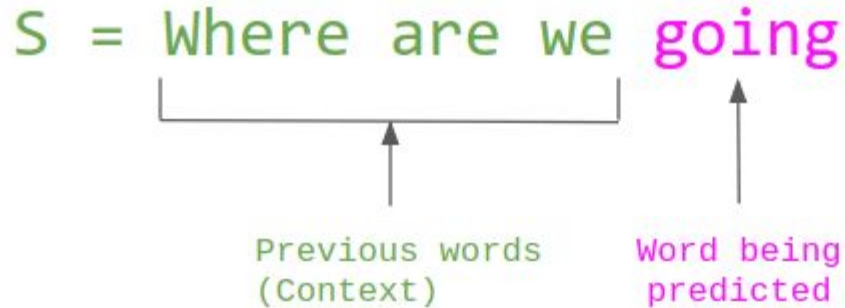
uses **LLMs** to generate tests based on provided bug reports ...

... then applies **post-processing** to ensure the developer only sees the best results.



Approach

Language Models are Autocomplete Machines



$$P(S) = P(\text{Where}) \times P(\text{are} \mid \text{Where}) \times P(\text{we} \mid \text{Where are}) \times P(\text{going} \mid \text{Where are we})$$

(image from *thegradients.pub*)

Formulating bug reproduction as autocomplete

Listing 1: Example prompt without examples.

```
1 # NaN in "equals" methods
2 ## Description
3 In "MathUtils", some "equals" methods will return true if both argument
   are NaN.
4 Unless I'm mistaken, this contradicts the IEEE standard.
5 If nobody objects, I'm going to make the changes. Report Content
6
7 ## Reproduction
8 >Provide a self-contained example that reproduces this issue.
9 ~ ~ ~
10 public void test
```

The first part of the prompt presents the bug report.

Formulating bug reproduction as autocomplete

Listing 1: Example prompt without examples.

```
1 # NaN in "equals" methods
2 ## Description
3 In "MathUtils", some "equals" methods will return true if both argument
   are NaN.
4 Unless I'm mistaken, this contradicts the IEEE standard.
5 If nobody objects, I'm going to make the changes.
6
7 ## Reproduction
8 >Provide a self-contained example that reproduces this issue.
9 ~ ~ ~
10 public void test
```

Prompting Reproducing Test Generation

The second part increases the likelihood of a bug-reproducing test
(from a language distribution perspective).

LLMs are known to benefit with examples

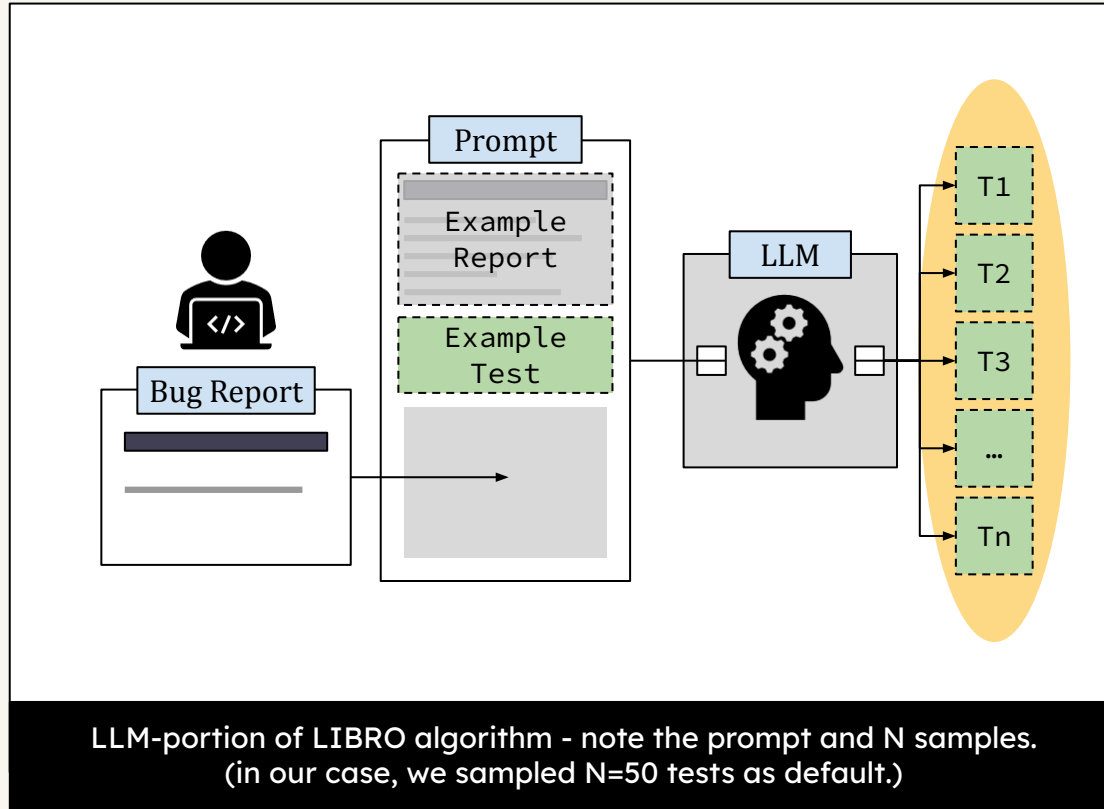
```

16 >Provide a self-contained example that reproduces this issue.
17 ...
18 public void testNumberUtils () {
19     assertEquals(Long.valueOf(0x80000000L),    NumberUtils.createNumber("0x80000000"));
20 }
21 ...
22
23 # #107 Incorrect date parsed when week and month used together
24 ## Description
25 I have following code snippet :
26
27 ```
28     DateTimeFormatter dtf = DateTimeFormat.forPattern("xxxxMM'w'ww");
29     DateTime dt = dtf.parseDateTime("201101w01");
30     System.out.println(dt);
31 ```
32
33 It should print 2011-01-03 but it is printing 2010-01-04.
34
35 Please let me know if I am doing something wrong here
36
37 ## Reproduction
38 >Provide a self-contained example that reproduces this issue.
39 ...
40 public void testIssue107() {
41     DateTimeFormatter dtf = DateTimeFormat.forPattern("xxxxMM'w'ww");
42     DateTime dt = dtf.parseDateTime("201101w01");
43     assertEquals(2011, dt.getYear());
44     assertEquals(1, dt.getMonthOfYear());
45     assertEquals(3, dt.getDayOfMonth());
46 }
47 ...
48
49 # {{title}}
50 ## Description
51 {{content}}
52 ## Reproduction
53 >Provide a self-contained example that reproduces this issue.
54 ...
55 public void test
56 {{endon}}:```
57

```

A prompt template we used for experiments.
Note the example answers (highlighted).

Given a prompt, sample N candidate tests.



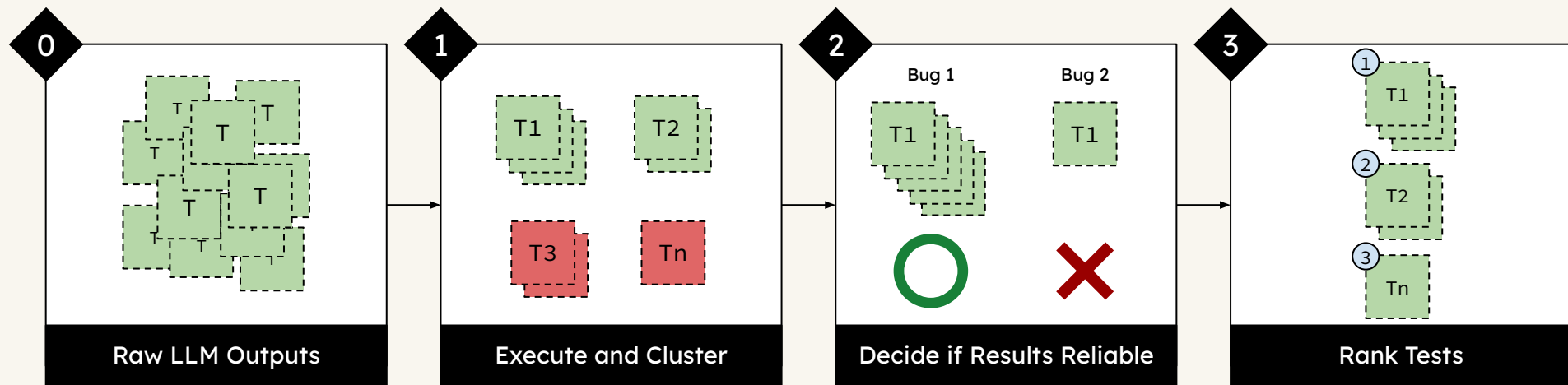
Showing 50 tests is infeasible

<pre>test1 { filler; filler2; }</pre>	<pre>test6 { filler; filler2; }</pre>	<pre>test11 { filler; filler2; }</pre>	<pre>test16 { filler; filler2; }</pre>	<pre>test21 { filler; filler2; }</pre>	<pre>test26 { filler; filler2; }</pre>	<pre>test31 { filler; filler2; }</pre>	<pre>test36 { filler; filler2; }</pre>	<pre>test41 { filler; filler2; }</pre>	<pre>test46 { filler; filler2; }</pre>
<pre>test2 { filler; filler2; }</pre>	<pre>test7 { filler; filler2; }</pre>	<pre>test12 { filler; filler2; }</pre>	<pre>test17 { filler; filler2; }</pre>	<pre>test22 { filler; filler2; }</pre>	<pre>test27 { filler; filler2; }</pre>	<pre>test32 { filler; filler2; }</pre>	<pre>test37 { filler; filler2; }</pre>	<pre>test42 { filler; filler2; }</pre>	<pre>test47 { filler; filler2; }</pre>
<pre>test3 { filler; filler2; }</pre>	<pre>test8 { filler; filler2; }</pre>	<pre>test13 { filler; filler2; }</pre>	<pre>test18 { filler; filler2; }</pre>	<pre>test23 { filler; filler2; }</pre>	<pre>test28 { filler; filler2; }</pre>	<pre>test33 { filler; filler2; }</pre>	<pre>test38 { filler; filler2; }</pre>	<pre>test43 { filler; filler2; }</pre>	<pre>test48 { filler; filler2; }</pre>
<pre>test4 { filler; filler2; }</pre>	<pre>test9 { filler; filler2; }</pre>	<pre>test14 { filler; filler2; }</pre>	<pre>test19 { filler; filler2; }</pre>	<pre>test24 { filler; filler2; }</pre>	<pre>test29 { filler; filler2; }</pre>	<pre>test34 { filler; filler2; }</pre>	<pre>test39 { filler; filler2; }</pre>	<pre>test44 { filler; filler2; }</pre>	<pre>test49 { filler; filler2; }</pre>
<pre>test5 { filler; filler2; }</pre>	<pre>test10 { filler; filler2; }</pre>	<pre>test15 { filler; filler2; }</pre>	<pre>test20 { filler; filler2; }</pre>	<pre>test25 { filler; filler2; }</pre>	<pre>test30 { filler; filler2; }</pre>	<pre>test35 { filler; filler2; }</pre>	<pre>test40 { filler; filler2; }</pre>	<pre>test45 { filler; filler2; }</pre>	<pre>test50 { filler; filler2; }</pre>

Some might not even compile!

<pre>test1 { filler; filler2; }</pre>	<pre>test6 { filler; filler2; }</pre>	<pre>test11 { filler; filler2; }</pre>	<pre>test16 { filler; filler2; }</pre>	<pre>test21 { filler; filler2; }</pre>	<pre>test26 { filler; filler2; }</pre>	<pre>test31 { filler; filler2; }</pre>	<pre>test36 { filler; filler2; }</pre>	<pre>test41 { filler; filler2; }</pre>	<pre>test46 { filler; filler2; }</pre>
<pre>test2 { filler; filler2; }</pre>	<pre>test7 { filler; filler2; }</pre>	<pre>test12 { filler; filler2; }</pre>	<pre>test17 { filler; filler2; }</pre>	<pre>test22 { filler; filler2; }</pre>	<pre>test27 { filler; filler2; }</pre>	<pre>test32 { filler; filler2; }</pre>	<pre>test37 { filler; filler2; }</pre>	<pre>test42 { filler; filler2; }</pre>	<pre>test47 { filler; filler2; }</pre>
<pre>test3 { filler; filler2; }</pre>	<pre>test8 { filler; filler2; }</pre>	<pre>test13 { filler; filler2; }</pre>	<pre>test18 { filler; filler2; }</pre>	<pre>test23 { filler; filler2; }</pre>	<pre>test28 { filler; filler2; }</pre>	<pre>test33 { filler; filler2; }</pre>	<pre>test38 { filler; filler2; }</pre>	<pre>test43 { filler; filler2; }</pre>	<pre>test48 { filler; filler2; }</pre>
<pre>test4 { filler; filler2; }</pre>	<pre>test9 { filler; filler2; }</pre>	<pre>test14 { filler; filler2; }</pre>	<pre>test19 { filler; filler2; }</pre>	<pre>test24 { filler; filler2; }</pre>	<pre>test29 { filler; filler2; }</pre>	<pre>test34 { filler; filler2; }</pre>	<pre>test39 { filler; filler2; }</pre>	<pre>test44 { filler; filler2; }</pre>	<pre>test49 { filler; filler2; }</pre>
<pre>test5 { filler; filler2; }</pre>	<pre>test10 { filler; filler2; }</pre>	<pre>test15 { filler; filler2; }</pre>	<pre>test20 { filler; filler2; }</pre>	<pre>test25 { filler; filler2; }</pre>	<pre>test30 { filler; filler2; }</pre>	<pre>test35 { filler; filler2; }</pre>	<pre>test40 { filler; filler2; }</pre>	<pre>test45 { filler; filler2; }</pre>	<pre>test50 { filler; filler2; }</pre>

LIBRO's post-processing in three steps

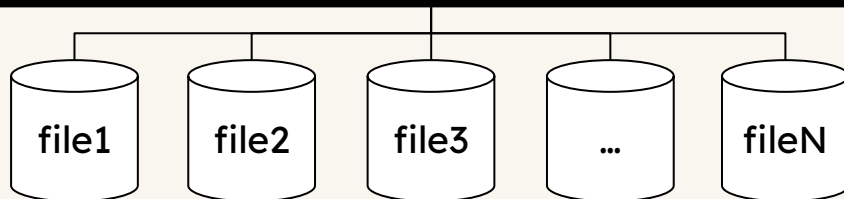


Injecting to target files

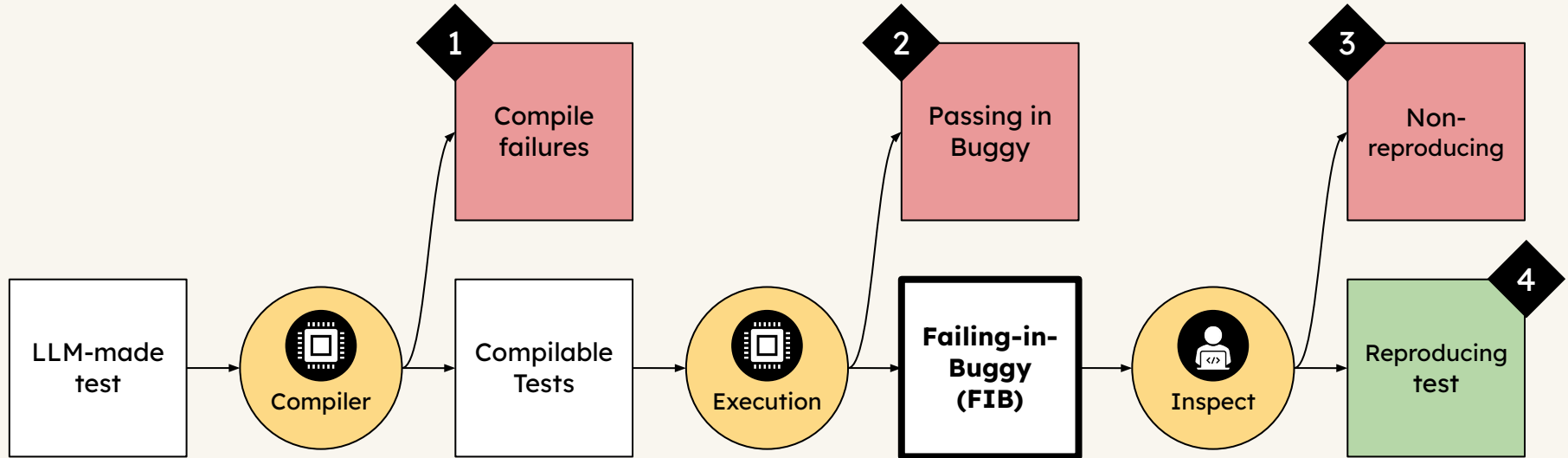
Listing 2: Example LLM result from the bug report described in Table II.

```
1 public void testEquals() {  
2     assertFalse(MathUtils.equals(Double.NaN, Double.NaN));  
3     assertFalse(MathUtils.equals(Float.NaN, Float.NaN));  
4 }
```

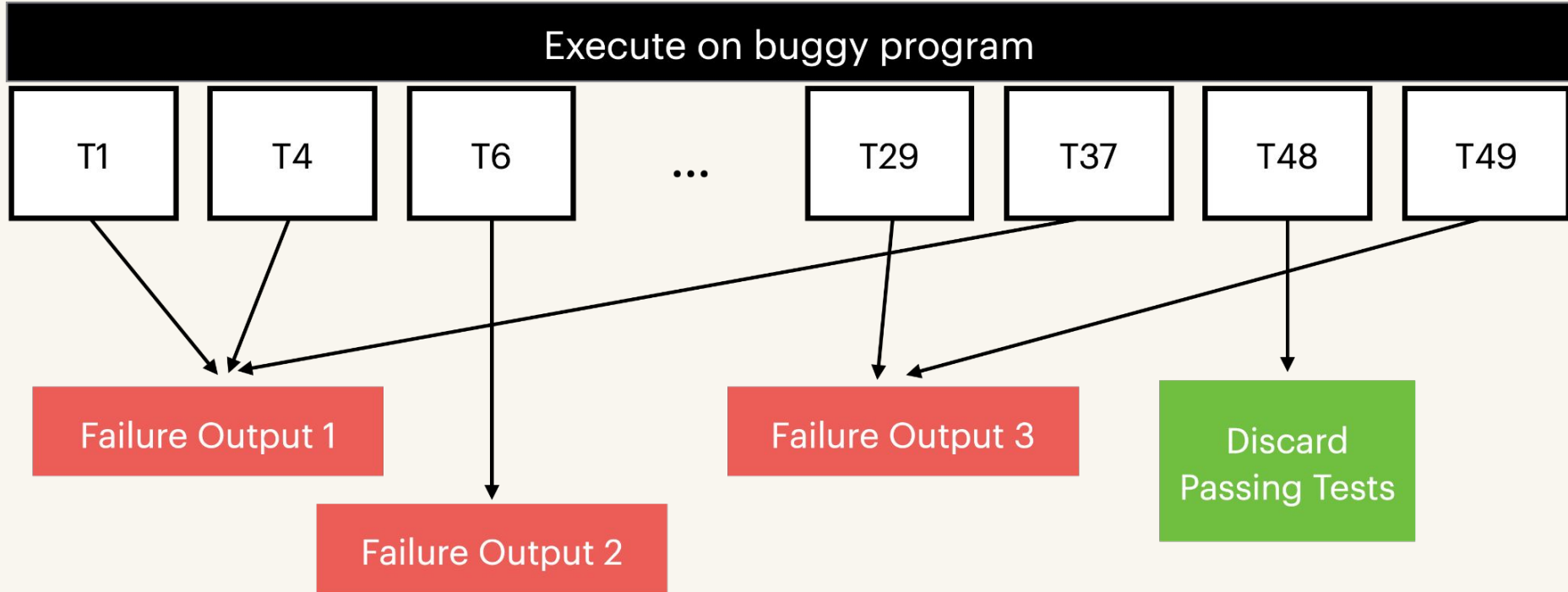
Select the file with greatest lexical similarity and inject the test; add `import` statements for unmet dependencies.



Execute Tests. Four results possible:

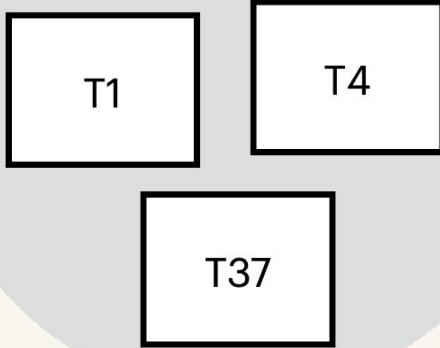


Cluster FIB tests with error message



Show results only if cluster size large enough

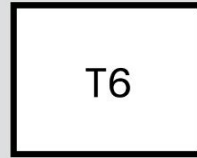
Failure Output 1



Cluster Size: **3**

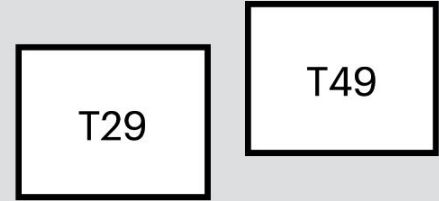
Accept the generated test set when
maximum cluster size > THRESHOLD

Failure Output 2



Cluster Size: 1

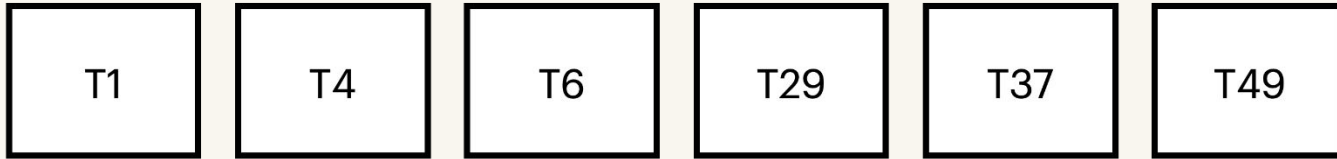
Failure Output 3



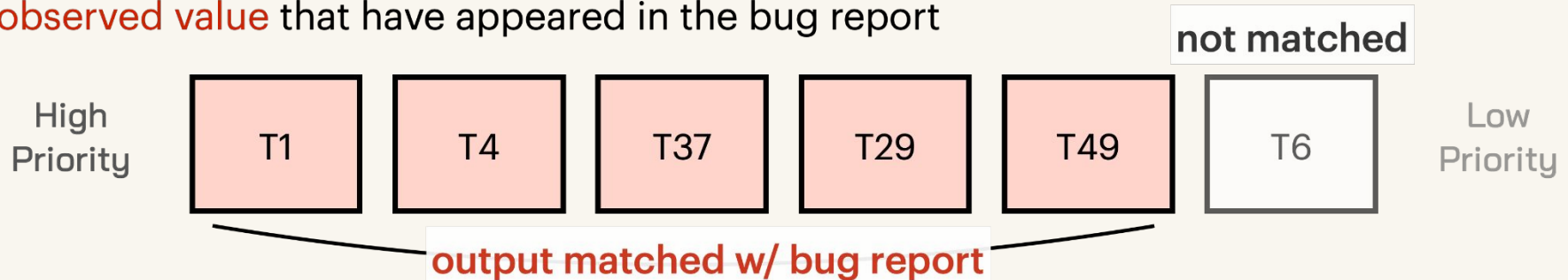
Cluster Size: 2

Ranking tests with three heuristics (1)

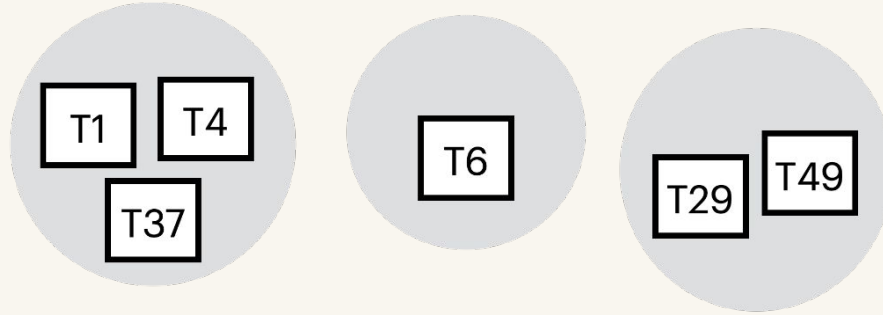
Q. Which test is more likely to be a correct bug reproducing test?



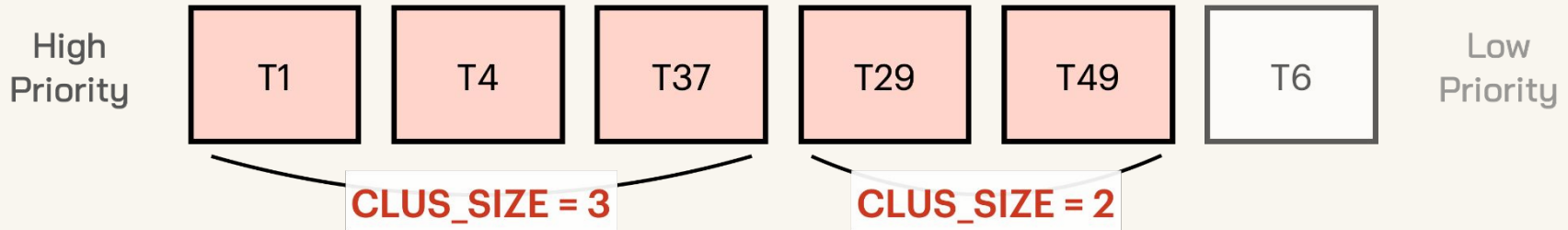
A1. **(Matching w/ Bug Report)** When test outputs **include exception type** or **observed value** that have appeared in the bug report



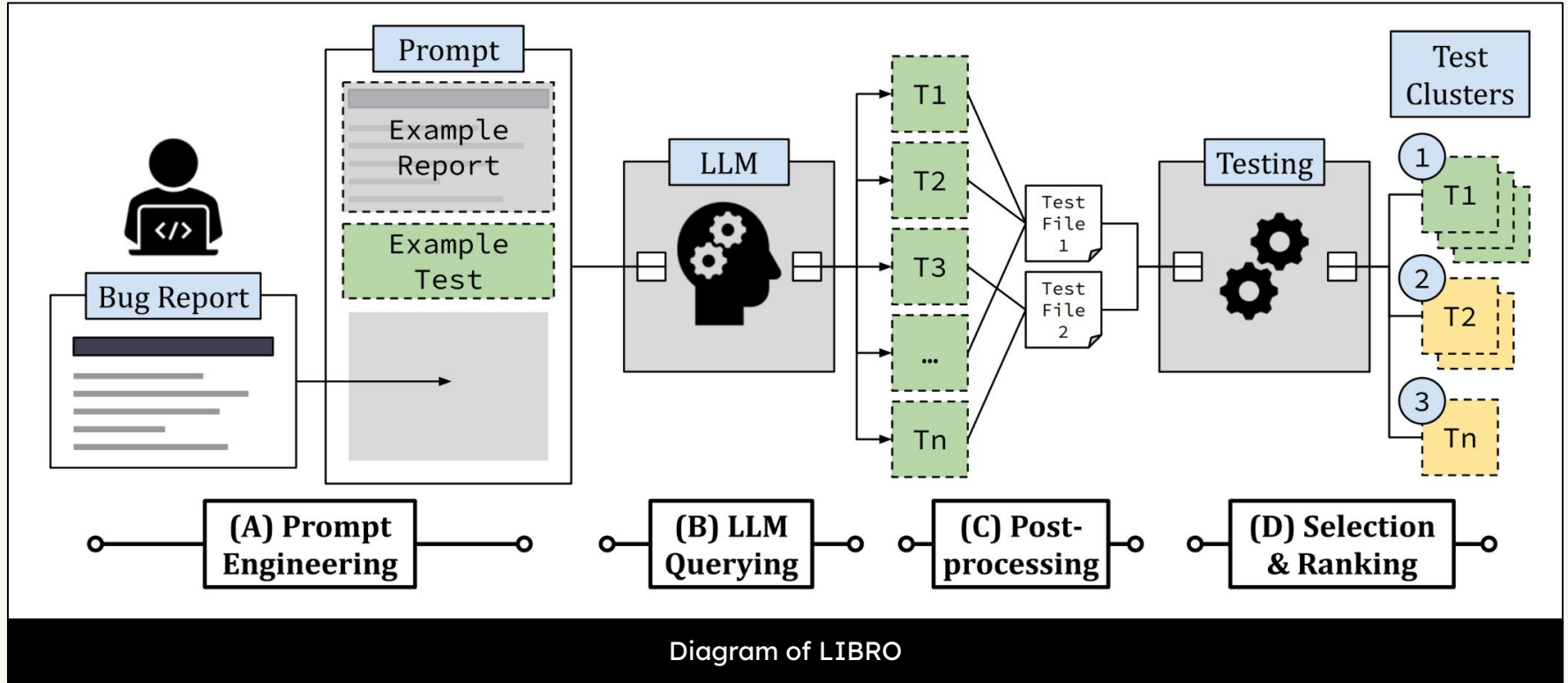
Ranking tests with three heuristics (2)



A2. (Consensus level) Tests from larger output cluster are prioritized



Recap

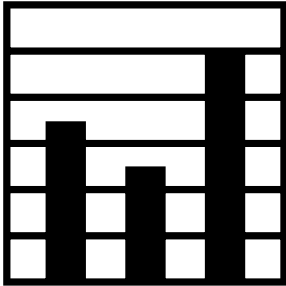




Results

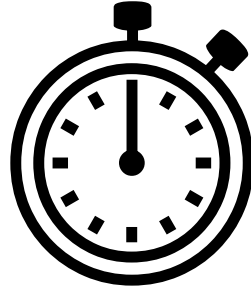
Evaluating the Technique

RQ1: Efficacy



How many bugs reproduced? (D4J)

RQ2: Efficiency



How much resources are required?

RQ3: Generality

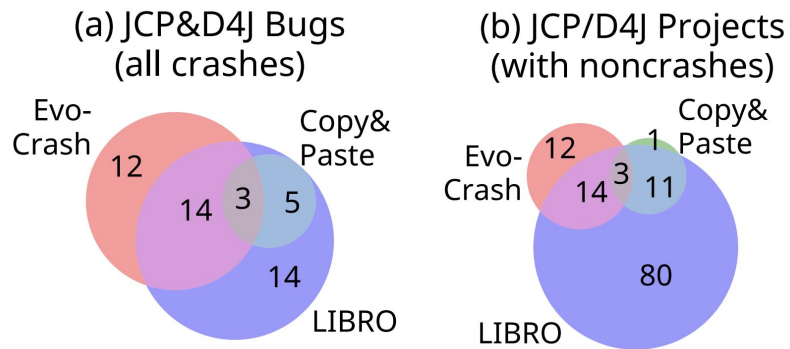


Are results similar on held-out data?

RQ1: Efficacy

Setting	reproduced	FIB
No Example (n=10)	124	440
One Example (n=10)	166	417
One Example from Source Project (n=10)	152	455
One Example with Constructor Info (n=10)	167	430
Two Examples (n=10, 5th percentile)	161	386
Two Examples (n=10, median)	173	409
Two Examples (n=10, 95th percentile)	184	429
Two Examples (n=50)	251	570
One Example, Crash Bugs (n=10)	69	153
One Example with Stack, Crash Bugs (n=10)	84	155

RQ1-1: One-third of all bugs were successfully reproduced.



RQ1-2: LIBRO significantly outperforms baselines.

RQ2-2: Time cost of each component

TABLE V: The time required for the pipeline of LIBRO

	Prompt	API	Processing	Running	Ranking	Total
Single Run	$<1 \mu\text{s}$	5.85s	1.23s	4.00s	-	11.1s
50-test Run	$<1 \mu\text{s}$	292s	34.8s	117s	0.02s	444s

The API call and actual execution of the test took the longest amount of time.

RQ2-3: Selection performance

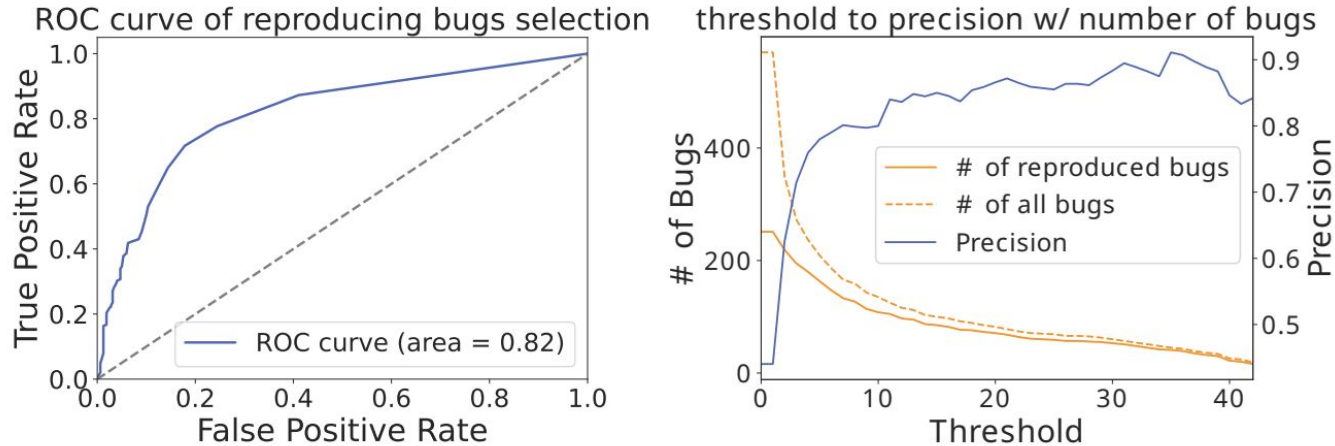


Fig. 4: ROC curve of bug selection (Left), Effect of thresholds to the number of bugs selected and precision (Right)

Selecting by counting the failing-in-buggy tests was effective;

RQ2-3: Ranking performance

TABLE VI: Ranking Performance Comparison between LIBRO and Random Baseline

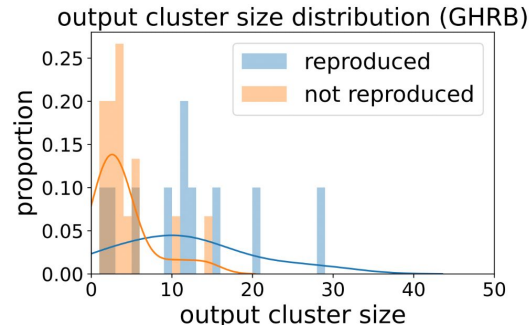
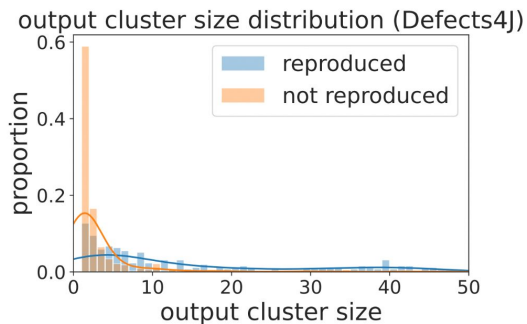
n	Defects4J				GHRB			
	$acc@n$ (precision)		$wef@n_{agg}$		$acc@n$ (precision)		$wef@n_{agg}$	
	LIBRO	random	LIBRO	random	LIBRO	random	LIBRO	random
1	149 (0.43)	116 (0.33)	201 (0.57)	234 (0.67)	6 (0.29)	4.8 (0.23)	15 (0.71)	16.2 (0.77)
3	184 (0.53)	172 (0.49)	539 (1.54)	599 (1.71)	7 (0.33)	6.6 (0.31)	42 (2.0)	44.6 (2.12)
5	199 (0.57)	192 (0.55)	797 (2.28)	874 (2.5)	8 (0.38)	7.3 (0.35)	60 (2.86)	64.3 (3.06)

Our within-FIB ranking technique was better than the random baseline.

RQ3: Generality

Project	rep/total	Project	rep/total	Project	rep/total
AssertJ	3/5	Jackson	0/2	Gson	4/7
checkstyle	0/13	Jsoup	2/2	sslcontext	1/2

RQ3-1: LIBRO demonstrating it works outside of its (potential) training data.



RQ3-2: LIBRO selection and ranking heuristics generalize to novel data as well.

Successful Example

DateTimeZone.forOffsetHoursMinutes cannot handle negative offset < 1 hour

`DateTimeZone.forOffsetHoursMinutes(h, m)` cannot handle negative offset < 1 hour like `-0:30` due to argument range checking. I used `forOffsetMillis()` instead.

This should probably be mentioned in the documentation or negative minutes be accepted.

Bug Report: Joda-Time BUG #8



```
public void testIssue() {
    DateTimeZone tz =
    DateTimeZone.forOffsetHoursMinutes(0, -30);
    assertNotNull(tz);
}
```

A test generated by LIBRO

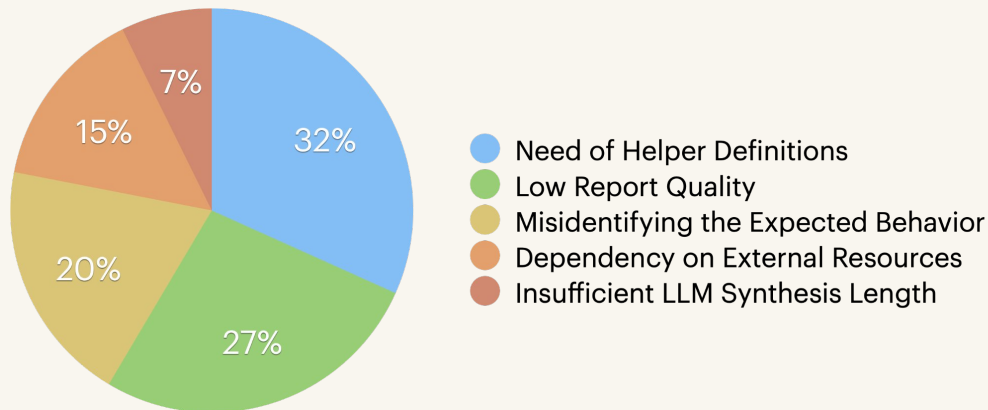


```
public void testForOffsetHoursMinutes_int_int() {
    ...
    assertEquals(DateTimeZone.forID("-00:15"),
    DateTimeZone.forOffsetHoursMinutes(0, -15));
    ...
}
```

A developer-written test

Examples and Failure Analysis

- When failures happened, we find that
 - 32.5% are due to a need of complex helper functions;
 - 27.5% are due to low report quality;
 - 20% are due to LLM misunderstanding of report;
 - 15% are due to dependency on external resources;
 - 7.5% are due to LLM synthesis limit (we set 256 tokens, or ~1000 characters).



Conclusion



1

We tackle the problem of **reproducing general bugs from reports.**

2

We propose LIBRO, which **combines LLMs and postprocessing** to effectively reproduce bug reports.

3

Our evaluation shows **LIBRO successfully reproduces bugs**, and that its postprocessing heuristics work.

Contact us at sungmin.kang@kaist.ac.kr / juyeon.yoon@kaist.ac.kr

Find our preprint with the QR code above, or by searching for “Exploring LLM-based General Bug Reproduction”



Extra Slides

Example 1 - Need of Helper Functions

- When failures happened, we find that
 - 32.5% are due to a need of complex helper functions;
 - 27.5% are due to low report quality;
 - 20% are due to LLM misunderstanding of report;
 - 15% are due to dependency on external resources;
 - 7.5% are due to LLM synthesis limit (we set 256 tokens, or ~1000 characters).

Invalid left-hand side of assignment not detected

What steps will reproduce the problem?

Compile this: `var x=0,y=1;x||y=8`

What is the expected/actual output?

I expect an error, because this is parsed as `(x||y)=8`, which is an invalid left-hand side of an assignment. Instead, I get `var x=0,y=1;x||y=8;`

Bug Report: Google Closure BUG #84



```
public void testIssue108() {
    try {
        ScriptEngine se = new
ScriptEngineManager().getEngineByName("JavaScript");
        se.eval("var x=0,y=1;x||y=8");
        fail("IllegalArgumentException was expected");
    } catch (ScriptException e) {
        // ok
    }
}
```

FAILS TO PERFECTLY REPRODUCE
THE SPECIALIZED TEST HELPER METHOD
(HERE, `parseError`)

A test generated by LIBRO(not bug-reproducing)

Example 1 - Need of Helper Functions

- When failures happened, we find that
 - 32.5% are due to a need of complex helper functions;
 - 27.5% are due to low report quality;
 - 20% are due to LLM misunderstanding of report;
 - 15% are due to dependency on external resources;
 - 7.5% are due to LLM synthesis limit (we set 256 tokens, or ~1000 characters).

Invalid left-hand side of assignment not detected

What steps will reproduce the problem?

Compile this: `var x=0,y=1;x||y=8`

What is the expected/actual output?

I expect an error, because this is parsed as `(x||y)=8`, which is an invalid left-hand side of an assignment. Instead, I get `var x=0,y=1;x||y=8;`

Bug Report: Google Closure BUG #84



```
/**
 * Verify that the given code has the given parse errors.
 * @return If in IDE mode, returns a partial tree.
 */
private Node parseError(String string, String... errors) {
    TestErrorReporter testErrorReporter = new TestErrorReporter(errors, null);
    Node script = null;
    try {
        StaticSourceFile file = new SimpleSourceFile("input", false);
        script = ParserRunner.parse(
            file, string, ParserRunner.createConfig(isIdeMode, mode, false),
            testErrorReporter, Logger.getAnonymousLogger()).ast;
    } catch (IOException e) {
        throw new RuntimeException(e);
    }
}
```


Example 2 - Low Report Quality

- When failures happened, we find that
 - 32.5% are due to a need of complex helper functions;
 - **27.5% are due to low report quality;**
 - 20% are due to LLM misunderstanding of report;
 - 15% are due to dependency on external resources;
 - 7.5% are due to LLM synthesis limit (we set 256 tokens, or ~1000 characters).

**StrBuilder contains usages of
thisBuf.length when they should use size**

While fixing LANG-294 I noticed that there are two other places in StrBuilder that reference thisBuf.length and unless I'm mistaken they shouldn't.

Bug Report: Defects4J Lang-60

External Dependency

- When failures happened, we find that
 - 32.5% are due to a need of complex helper functions;
 - 27.5% are due to low report quality;
 - 20% are due to LLM misunderstanding of report;
 - **15% are due to dependency on external resources;**
 - 7.5% are due to LLM synthesis limit (we set 256 tokens, or ~1000 characters).

PAX header entry name ending with / causes problems

There seems to be a problem when a PAX header entry (link flag is 'x') has a name ending with "/". The `TarArchiveEntry.isDirectory()` check ends up returning true because of the trailing slash which means no content can be read from the entry (...)

Bug Report: Commons-Compress BUG #38



```
public void testPaxHeaderFails() throws Exception {
    FileInputStream fin = null;
    try {
        fin = new
        FileInputStream("/path/to/docker-image-layer.tar");
        TarArchive tar = new TarArchive(fin);
        TarArchiveInputStream(fin);
        (...)
    }
}
```

**Fails due to missing file;
The generated test is dependent
on an external file**

A test generated by LIBRO (not bug-reproducing)

Synthesis limit example

- When failures happened, we find that
 - 32.5% are due to a need of complex helper functions;
 - 27.5% are due to low report quality;
 - 20% are due to LLM misunderstanding of report;
 - 15% are due to dependency on external resources;
 - **7.5% are due to LLM synthesis limit (we set 256 tokens, or ~1000 characters).**

```
* unparsable
...
public void testNumberUtilsCreateNumber() {
    assertEquals(Long.valueOf(1L), NumberUtils.createNumber("1l"));
    assertEquals(Long.valueOf(1L), NumberUtils.createNumber("1L"));
    assertEquals(Long.valueOf(1L), NumberUtils.createNumber("+1l"));
    assertEquals(Long.valueOf(1L), NumberUtils.createNumber("+1L"));
    assertEquals(Long.valueOf(-1L), NumberUtils.createNumber("-1l"));
    assertEquals(Long.valueOf(-1L), NumberUtils.createNumber("-1L"));
    assertEquals(Long.valueOf(1L), NumberUtils.createNumber("1"));
    assertEquals(Long.valueOf(1L), NumberUtils.createNumber("L"));
    assertEquals(Long.valueOf(1L), NumberUtils.createNumber("+1"));
    assertEquals(Long.valueOf(1L), NumberUtils.createNumber("
...

```