

When to Let the Developer Guide: Trade-offs between Open and Guided Test Amplification

Carolin Brandt
Danyao Wang
Andy Zaidman



Test Amplification

Generate x Unit tests by mutating existing tests

Existing Test →

→ Amplified Tests

@Test

```
void text_escape () {  
    A a = new A ("text");  
    B b = a.escape();  
    assertEquals (b.toString(), "text");  
}
```



Test Amplification

Generate xUnit tests by mutating existing tests



@Test

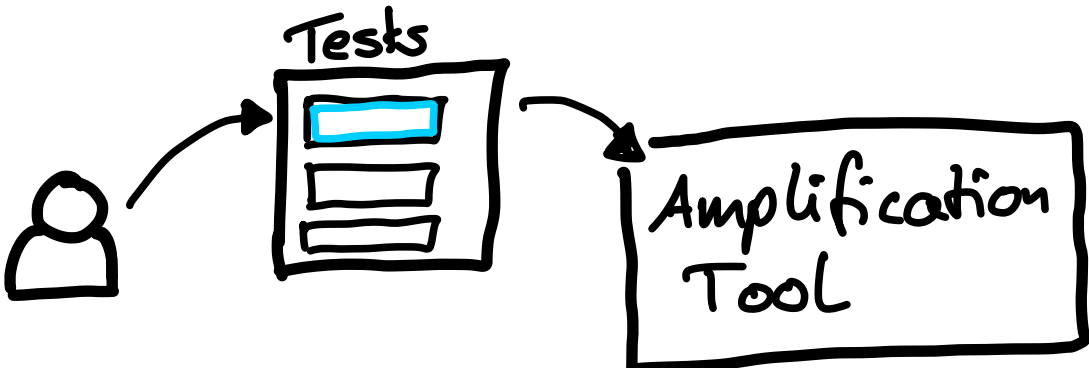
```
void text_escape () {  
    A a = new A ("text"); → "text $"  
    B b = a.escape()  
    assertEquals (b.toString(), "text"); → "text \\$\\ \\ \\"  
}
```

Open Test Amplification

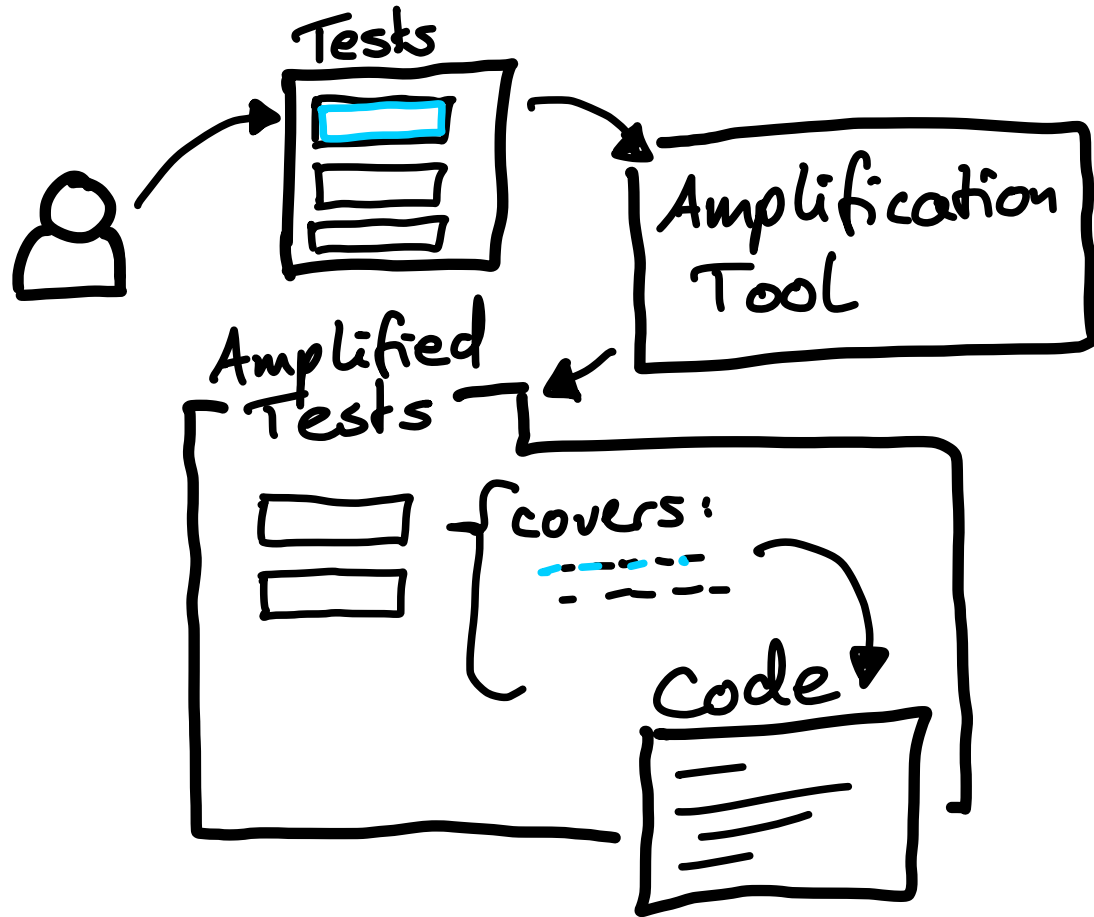
Looking in all directions for tests



User Interaction

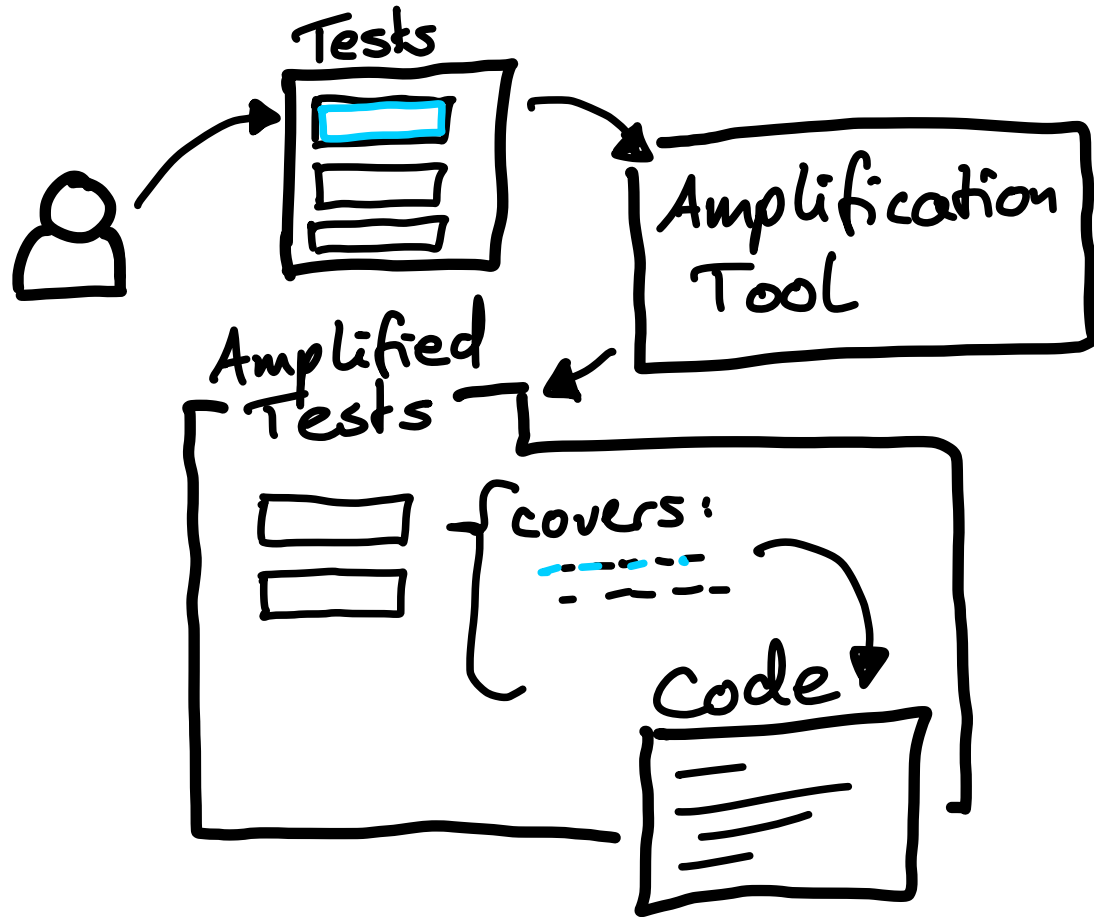


User Interaction



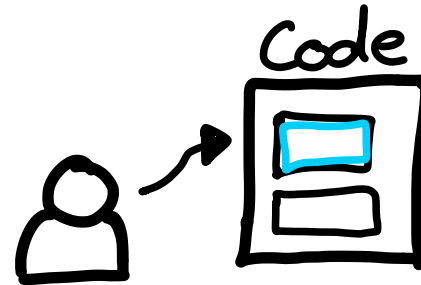
Open Test Amplification:
Looking in all Directions for Tests

User Interaction

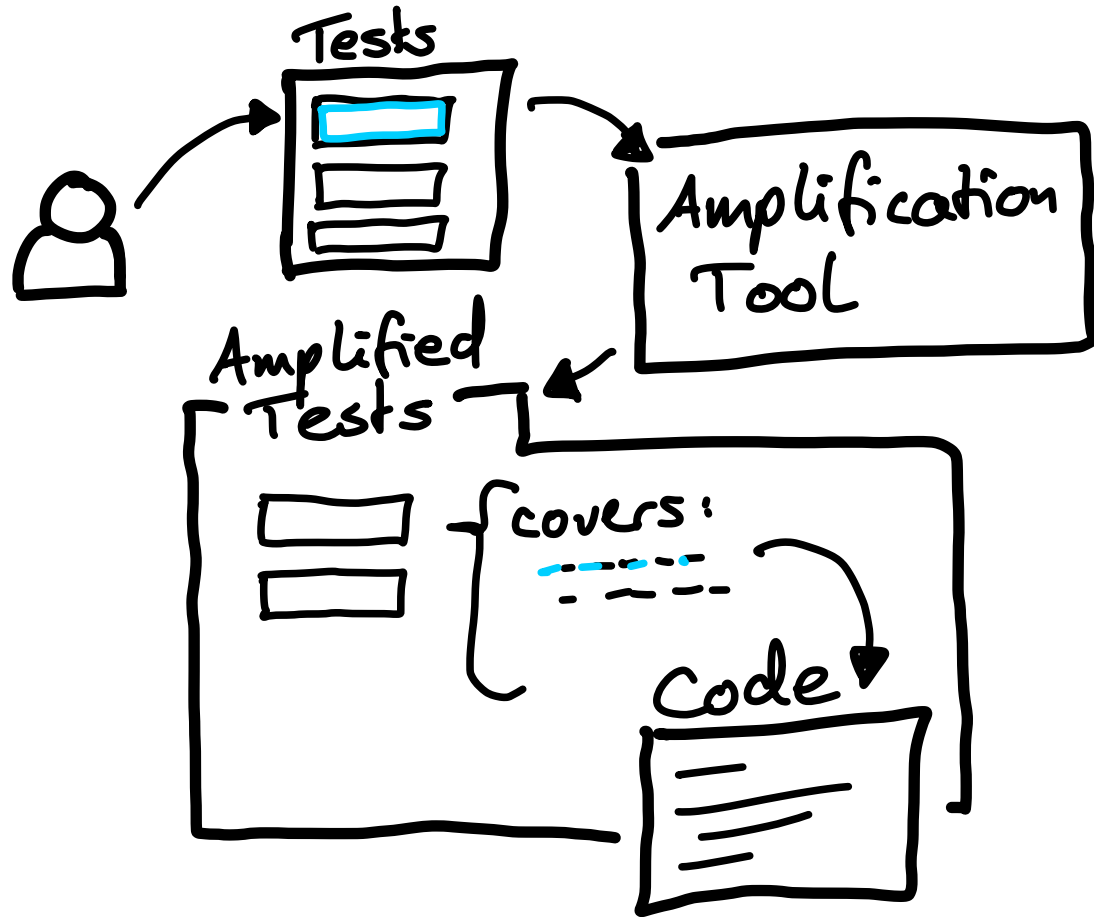


Open Test Amplification:
Looking in all Directions for Tests

What if we let the developer point us to what they want to test?

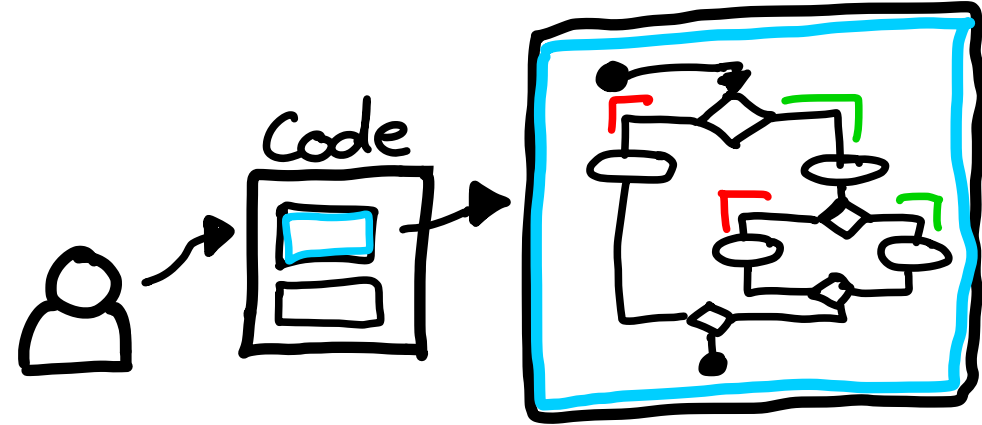


User Interaction

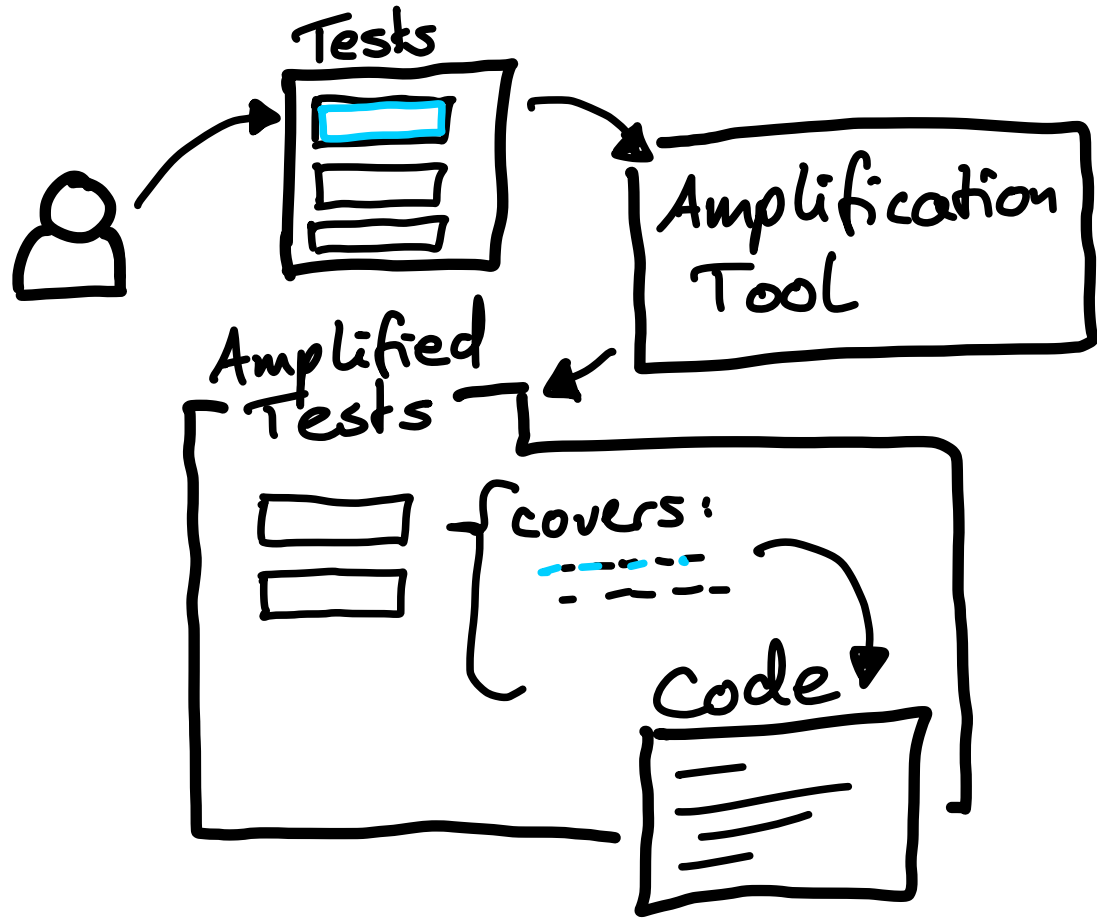


Open Test Amplification:
Looking in all Directions for Tests

What if we let the developer point us to what they want to test?

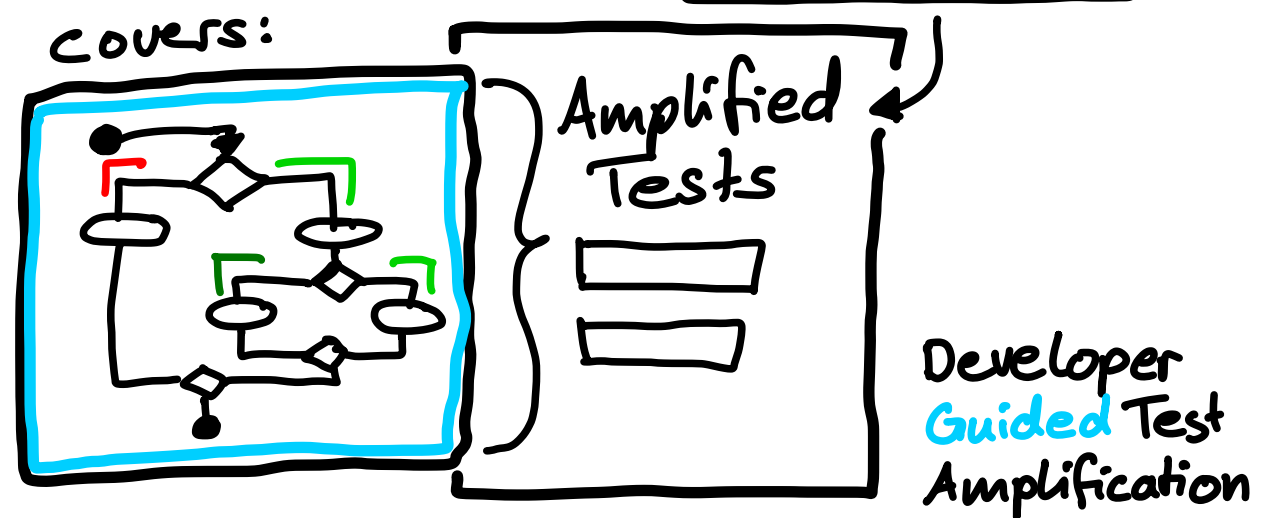
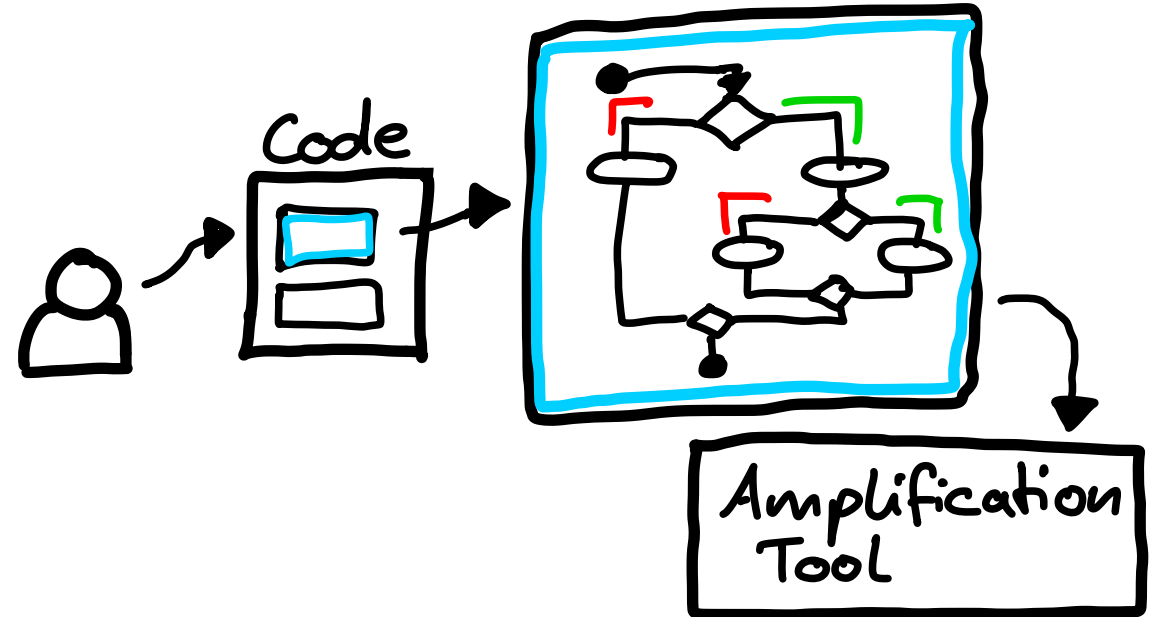


User Interaction

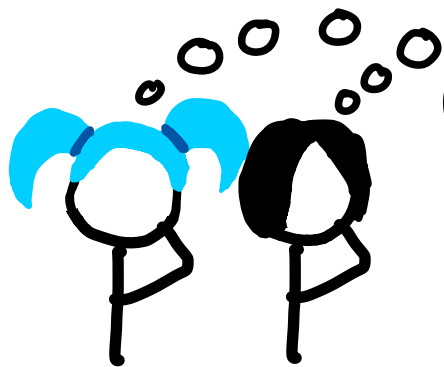


Open Test Amplification:
Looking in all Directions for Tests

What if we let the developer point us to what they want to test?

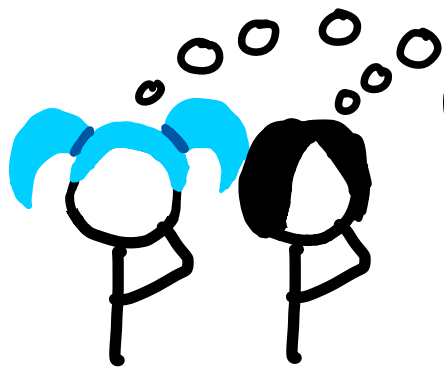
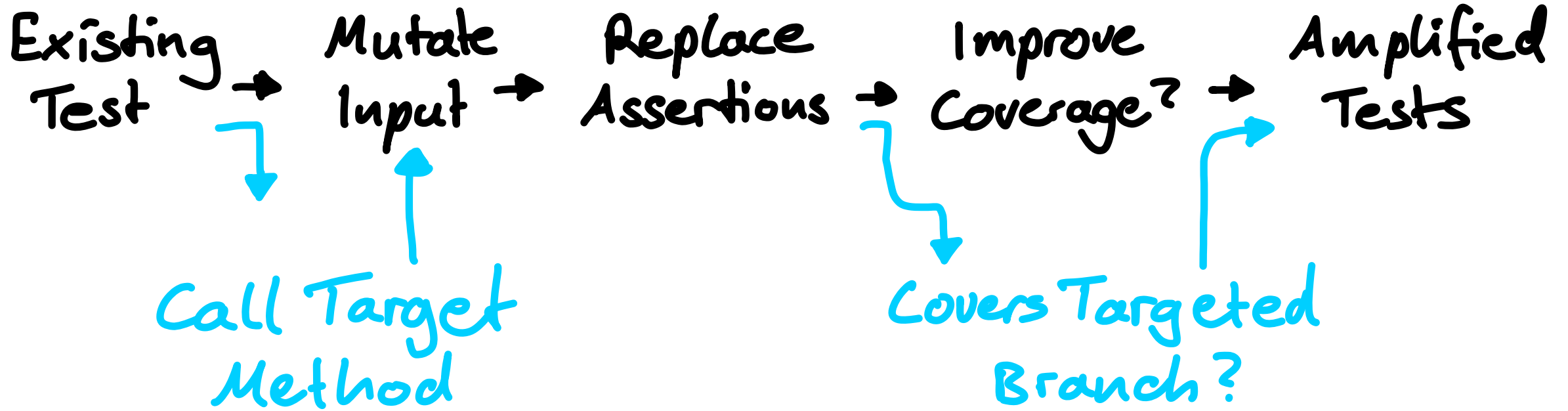


Guided Test Amplification



How effective are simple modifications?

Guided Test Amplification



How effective are simple modifications?

Technical Evaluation

sample 100 branches each from 2 projects

Technical Evaluation

sample 100 branches each from 2 projects

Sampled branches covered in total

	Javapoet	Stream-lib
Open	23%	35%
Guided	32%	41%

Main culprit:

correct initialization of objects under test

to hit targeted branch

Technical Evaluation

sample 100 branches each from 2 projects

Sampled branches covered in total
Javapoet Stream-lib

Open	23%	35%
Guided	32%	41%

Tests covering target branch per run
Javapoet Stream-lib

Open	24%	45%
Guided	70%	70%

Main culprit:

correct initialization of objects under test

to hit targeted branch

User Study Evaluation

Generate Tests with Open + Guided Test Amplification with IntelliJ Plugin

12 participants in 4 groups, 2 classes, crossover

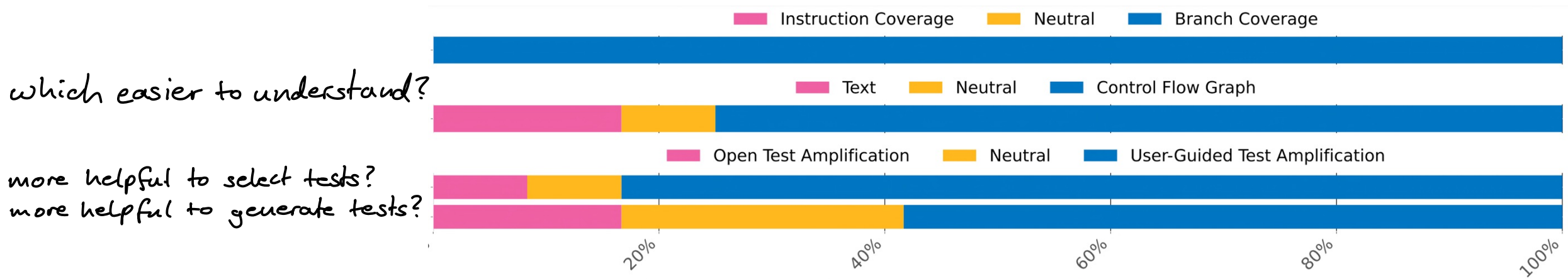
Semi-structured Interviews

User Study Evaluation

Generate Tests with Open + Guided Test Amplification with IntelliJ Plugin

12 participants in 4 groups, 2 classes, crossover

Semi-structured Interviews



Guided is not necessarily better...

Guided is not necessarily better...

Trade-offs: Guided vs. Open Test Amplification

Guided is not necessarily better...

Trade-offs: Guided vs. Open Test Amplification

Fits use case	Writing code & wanting tests for it	Improving test suite
---------------	-------------------------------------	----------------------

Guided is not necessarily better...

Trade-offs: Guided vs. Open Test Amplification

Fits use case	Writing code & wanting tests for it	Improving test suite
Understand Coverage	In target method in detail	Across whole project

Guided is not necessarily better...

Trade-offs: Guided vs. Open Test Amplification

Fits use case	Writing code & wanting tests for it	Improving test suite
Understand Coverage	In target method in detail	Across whole project
Expectation of receiving tests	Might disappoint if target branch cannot be covered	Only proposes tests it can provide

Guided is not necessarily better..


Trade-offs: Guided vs. Open Test Amplification

Fits use case	Writing code & wanting tests for it	Improving test suite
Understand Coverage	In target method in detail	Across whole project
Expectation of receiving tests	Might disappoint if target branch cannot be covered	Only proposes tests it can provide
Runtime efficiency	More effective at providing tests for targeted method	Larger variety of tests for whole class

Trade-offs: Guided vs. Open Test Amplification

Fits use case	Writing code & wanting tests for it	Improving test suite
Understand coverage	In target method in detail	Across whole project
Expectation of receiving tests	Might disappoint if target branch cannot be covered	Only proposes tests it can provide
Runtime efficiency	More effective at providing tests for targeted method	Larger variety of tests for whole class

When to Let the Developer Guide: Trade-offs Between Open and Guided Test Amplification

Carolin Brandt, Danyao Wang, Andy Zaidman 



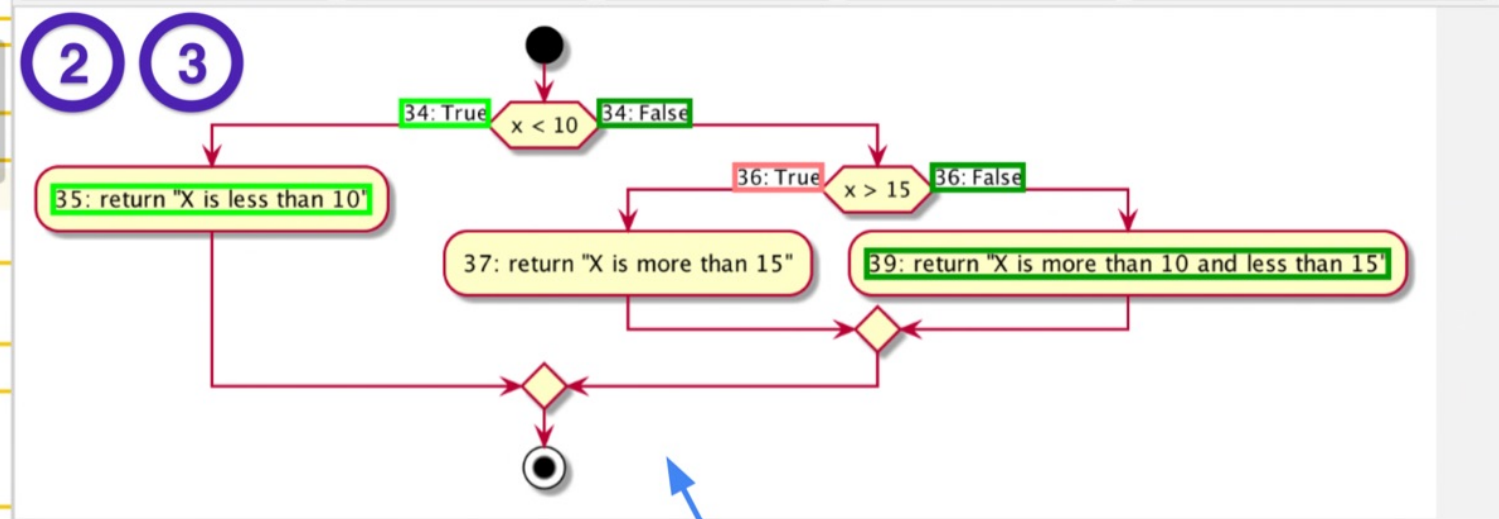
```
Calculator.java x CalculatorTest.java x
32
33 1 public String range() {
34     if (x < 10) {
35         return "X is less than 10";
36     } else if (x > 15) {
37         return "X is more than 15";
38     } else {
39         return "X is more than 10 and less than 15";
40     }
41 }
42
43 public String SingleIf(){
44     if (y > x) {
45         return "Y is larger than X";
46     }
47     return "X is larger than Y";
48 }
49
50 public int uncovered() {
51     for (int i = 1; i < x; ++i) {
52         --x;
53     }

```

```
Test Cube Test Generation for range()'
1 @Test(timeout = 10000)
2 public void testSub_mg12_assSep300() throws Exception {
3     Assert.assertEquals("X is less than 10", new Calculator(3, 2).other());
4 }

```

Add Test To Test Suite Ignore Test Case Next Test Case Previous Test Case Close Amplification Result



Test Cube found 3 amplified test cases.

Actions

- Inspect amplification results
- Inspect DSpot terminal output

Generate test to cover the selected branch Close

blems Terminal Build Dependencies

3 Event Log

```
AttributeTest.java x
1 package org.jsoup.nodes;
2
3 import ...
8
9 public class AttributeTest {
10     @Test
11     public void html() {
12         Attribute attr = new Attribute("key", "value &");
13         assertEquals("expected: \"key=\\\"value &amp;\\\"", attr.html());
14         assertEquals(attr.html(), attr.toString());
15     }
16
17     @Test public void testWithSupplementaryCharacterInAttributeKeyAndValue() {
18         String s = new String(Character.toChars(0x1F600));
19         Attribute attr = new Attribute(s, "A" + s + "B");
20         assertEquals("expected: \"A" + s + "B\"", attr.html());
21         assertEquals(attr.html(), attr.toString());
22
23     }
24     @Test public void validatesKeysNotEmpty() {
25         assertThrows(IllegalArgumentException.class, () -> new Attribute(" ", "value"));
26     }
27
28     @Test public void validatesKeysNotEmptyViaSet() {
29         assertThrows(IllegalArgumentException.class, () -> {
30             Attribute attr = new Attribute("One", "Check");
31             attr.setKey(" ");
32         });
33 }
```

Test Cube Amplification of 'html()'

Amplified test case 'html_assSep5'

Input modifications: 0 **4**
Assert statements added: 1

This test case improves the coverage in these classes/methods/lines:
(Click on the green links to see these lines within the class)

[org.jsoup.nodes.Attribute.hashCode](#)
L. 198 +8 instr.
L. 199 +12 instr.
L. 200 +2 instr.

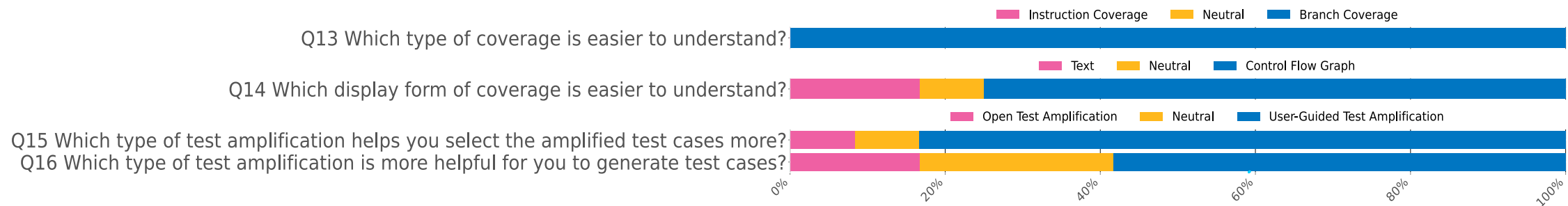
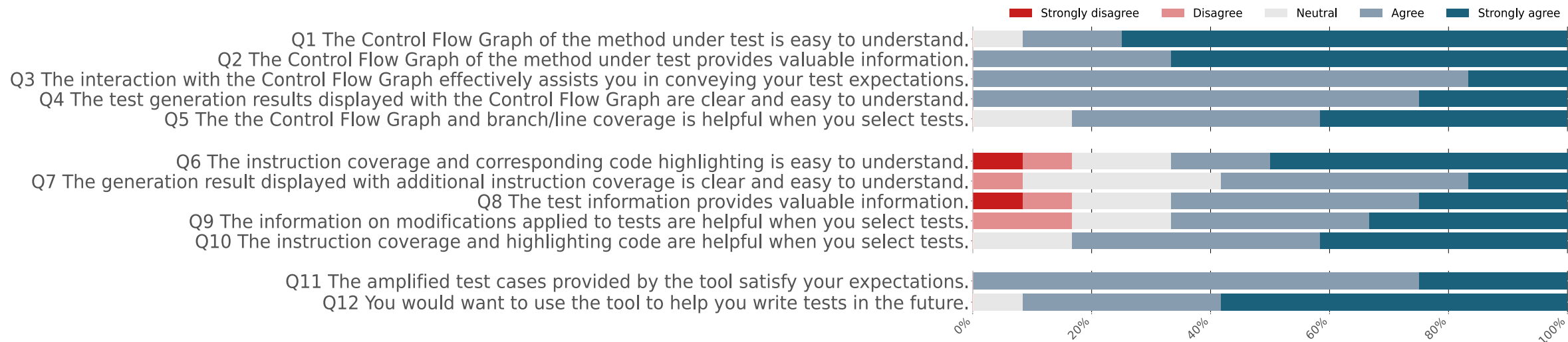
6 Add Test To Test Suite Ignore Test Case Next Test Case Previous Test Case Close

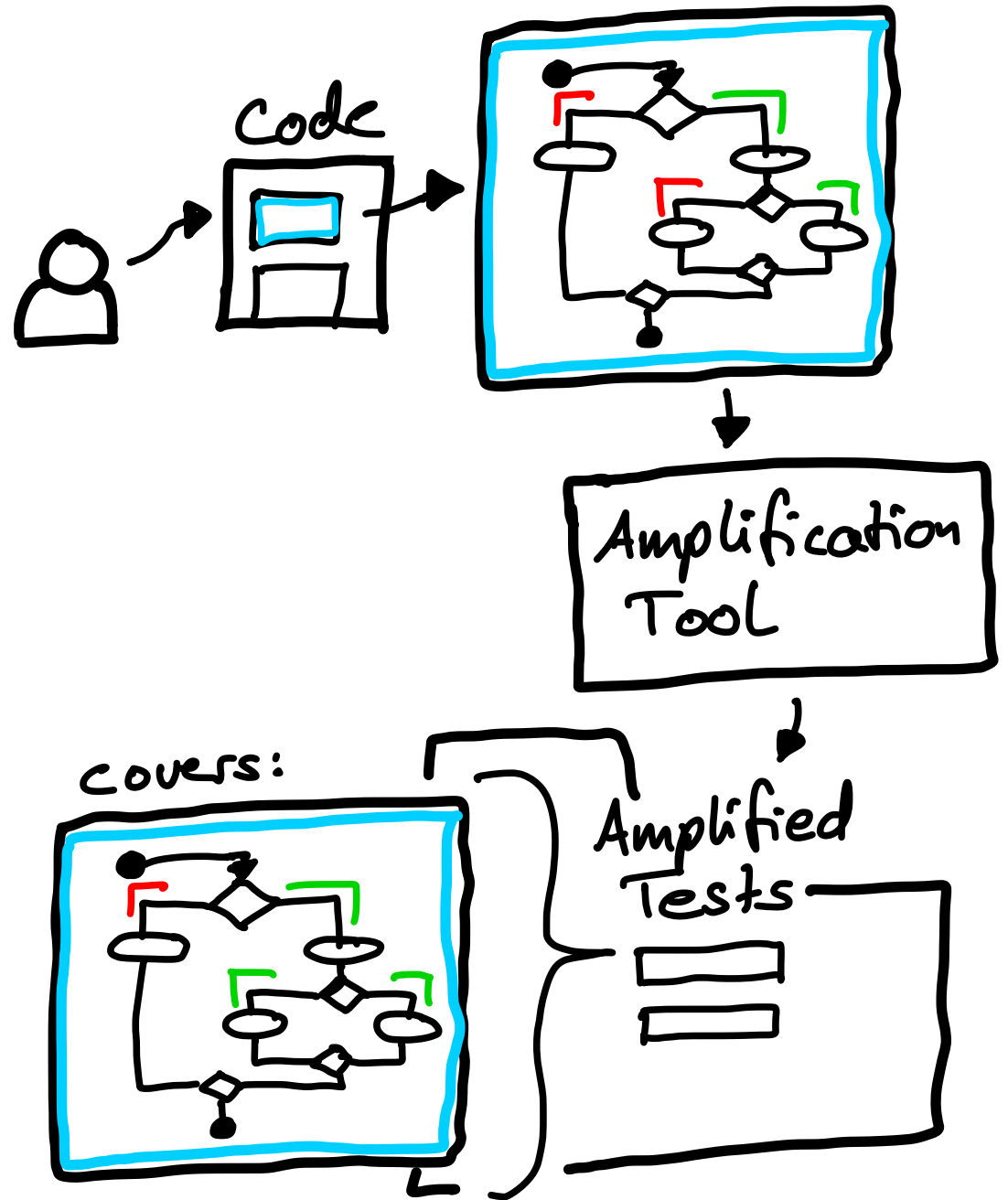
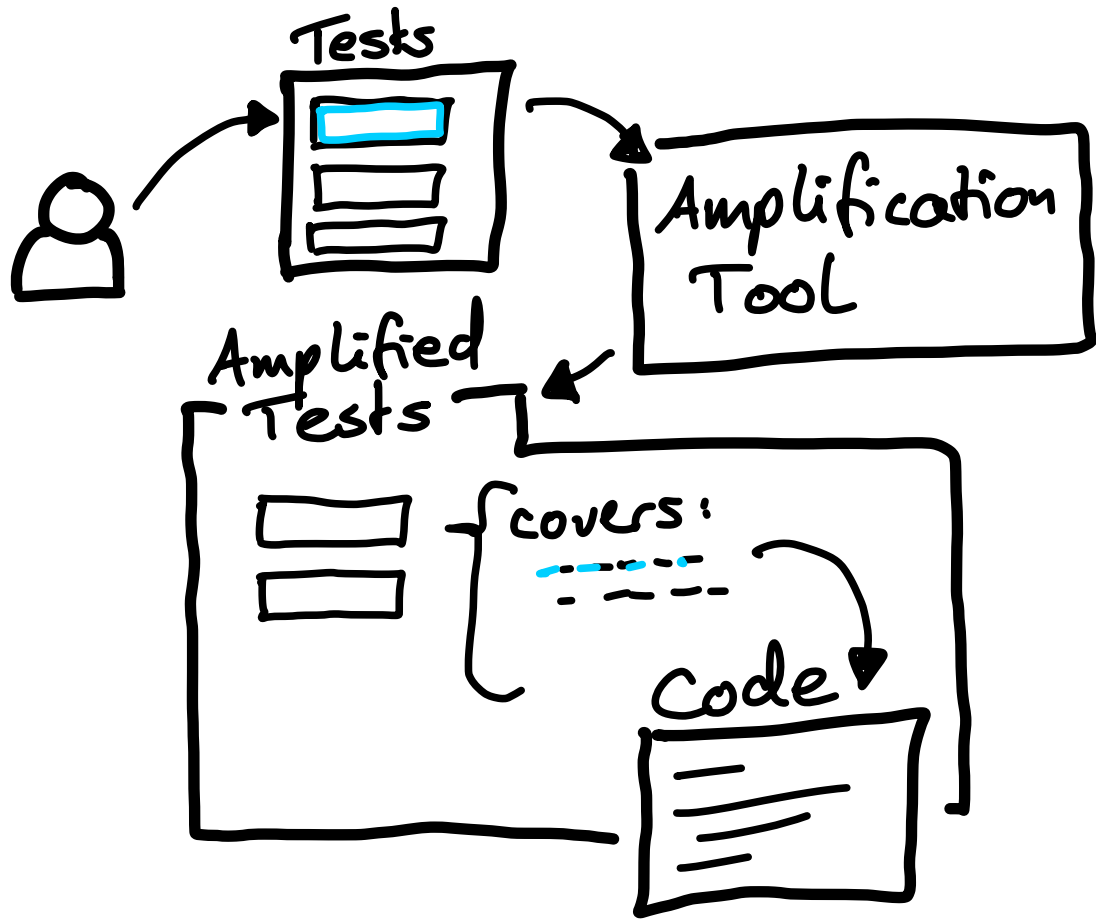
```
6 public void html_assSep5() throws Exception {
7     Attribute attr = new Attribute("key", "value &");
8     Assertions.assertEquals(234891960, ((int) ((Attribute) (attr)).hashCode()));
9 }
10
11 @Test
12 public void html_literalMutationString19_assSep92() throws Exception {
13     Attribute attr = new Attribute("key", " Hello\nthere NBSP ");
14     Assertions.assertEquals("key=\\\" Hello\nthere &nbsp; \\\", ((Attribute) (attr)).toString());
15 }
16
17 @Test
18 public void html_mg22_assSep1
```

3 Test Cube found 4 amplified test cases.
[Inspect amplification results](#) [Inspect DSpot terminal o](#)

2 Amplify 'testWithSupplementaryCharacterInAttributeKeyAndValue()'

3 Test Cube found 4 amplified test cases.
[Inspect amplification results](#) [Inspect DSpot terminal o](#)





1) To be actually helpful for users, we need to consciously design how users interact with generative tools like GPT, Copilot, ...

2) Too much research blindly applies LLMs to inappropriate tasks (tasks that require other skills than (natural) language composition).