# CS 563:
# Software Maintenance and Evolution

# **Automated Program Repair**

### Oregon State University, Spring 2024

# Motivation

**606   SOFTWARE FAILURES**

**1177   NEWS STORIES**

**314   COMPANIES**

LO$$E$ FROM SOFTWARE FAILURES (USD)

# 1,715,430,778,504

ONETRILLIONSEVENHUNDREDFIFTEENBILLIONFOURHUNDREDTHIRTYMILLIONSEVENHUNDREDSEVENTY-EIGHTTHOUSANDFIVEHUNDREDFOUR

3.7 OF 7.4
World Population
in Billions

PEOPLE AFFECTED (AT LEAST)
3,683,212,665

ACCUMULATED TIME     LOST
46 MINUTES
8 HOURS
3 DAYS
2 WEEKS
8 MONTHS
268 YEARS

Source: Software Fail Watch: 5th Edition, 2018.
Available at: https://www.tricentis.com/software-fail-watch

*"Legacy software testing tools that were developed two decades ago—some which are still heavily in use today—were never intended to support high levels of quality in today's rapid release cycles. It's time for a change."*

# Automatic Program Repair

Department: xxxx
Editor: Name, xxxx@email

## On the Introduction of Automatic Program Repair in Bloomberg

Serkan Kirbas, Etienne Windels, Olayori McBello, Kevin Kells, Matthew Pagano, Rafal Szalanski,
Bloomberg, London, UK & New York, USA

Vesna Nowack[1], Emily Winter[2], Steve Counsell[3], David Bowes[2], Tracy Hall[2], Saemundur Haraldsson[4], John Woodward[1]
[1]Queen Mary, University of London, UK
[2]Lancaster University, UK
[3]Brunel University, London, UK
[4]University of Stirling, UK

*Abstract*—A key to the success of Automatic Program Repair techniques is how easily they can be used in an industrial setting. In this article, we describe a collaboration by a team from four UK-based universities with Bloomberg (London) in implementing automatic, high-quality fixes to its code base. We explain the motivation for adopting APR, the mechanics of the prototype tool that was built, and the practicalities of integrating APR into existing systems.
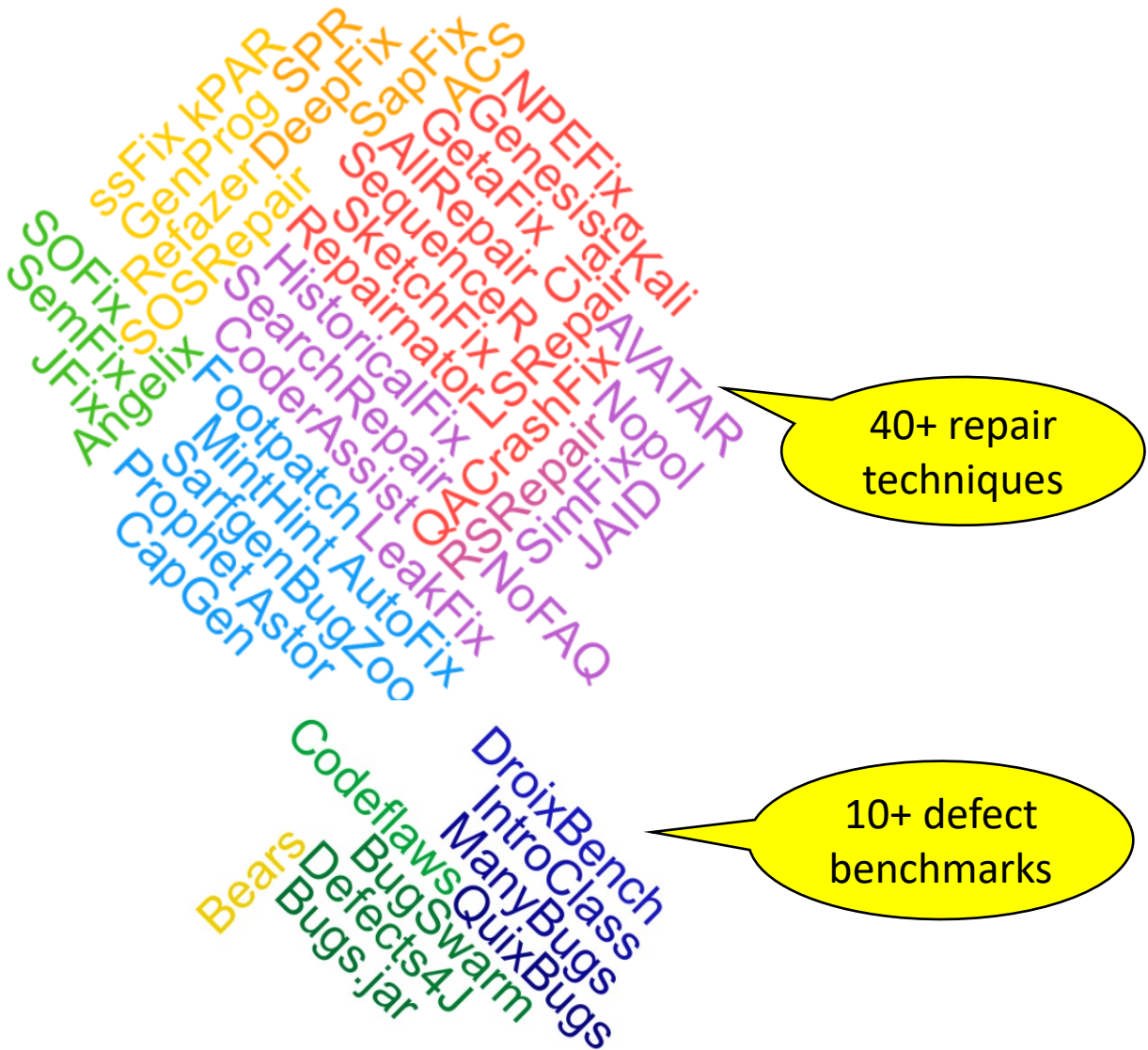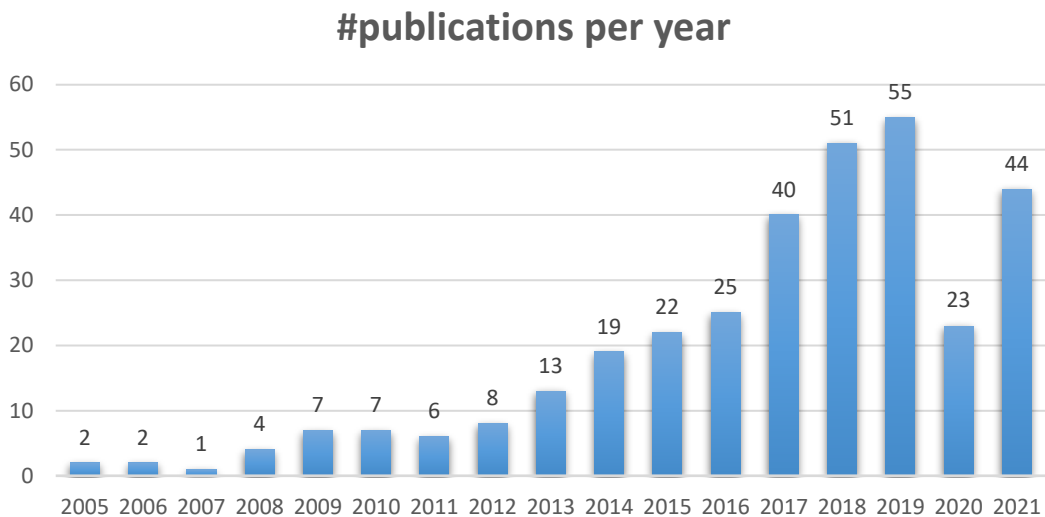
**facebook** Engineering

| Open Source ⌄ | Platforms ⌄ | Infrastructure Systems ⌄ | Physical Infrastructure ⌄ | Video Engineering & AR/VR ⌄ |

POSTED ON SEP 13, 2018 TO AI RESEARCH, DEVELOPER TOOLS, OPEN SOURCE, PRODUCTION ENGINEERING

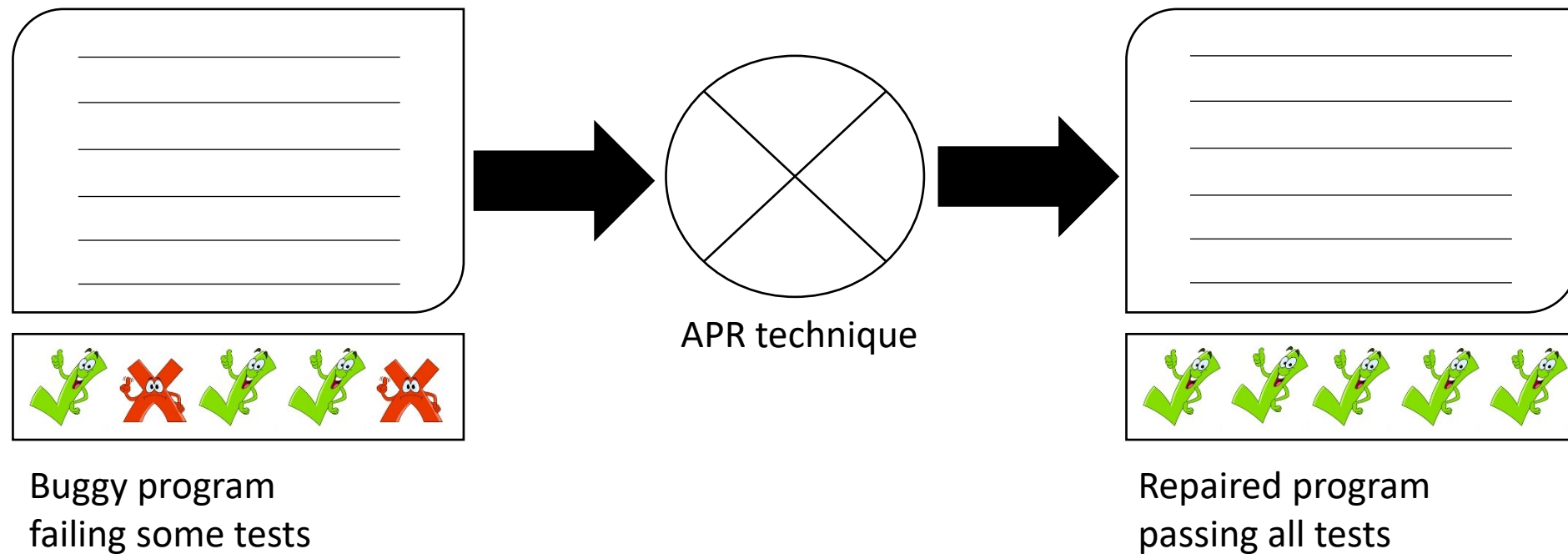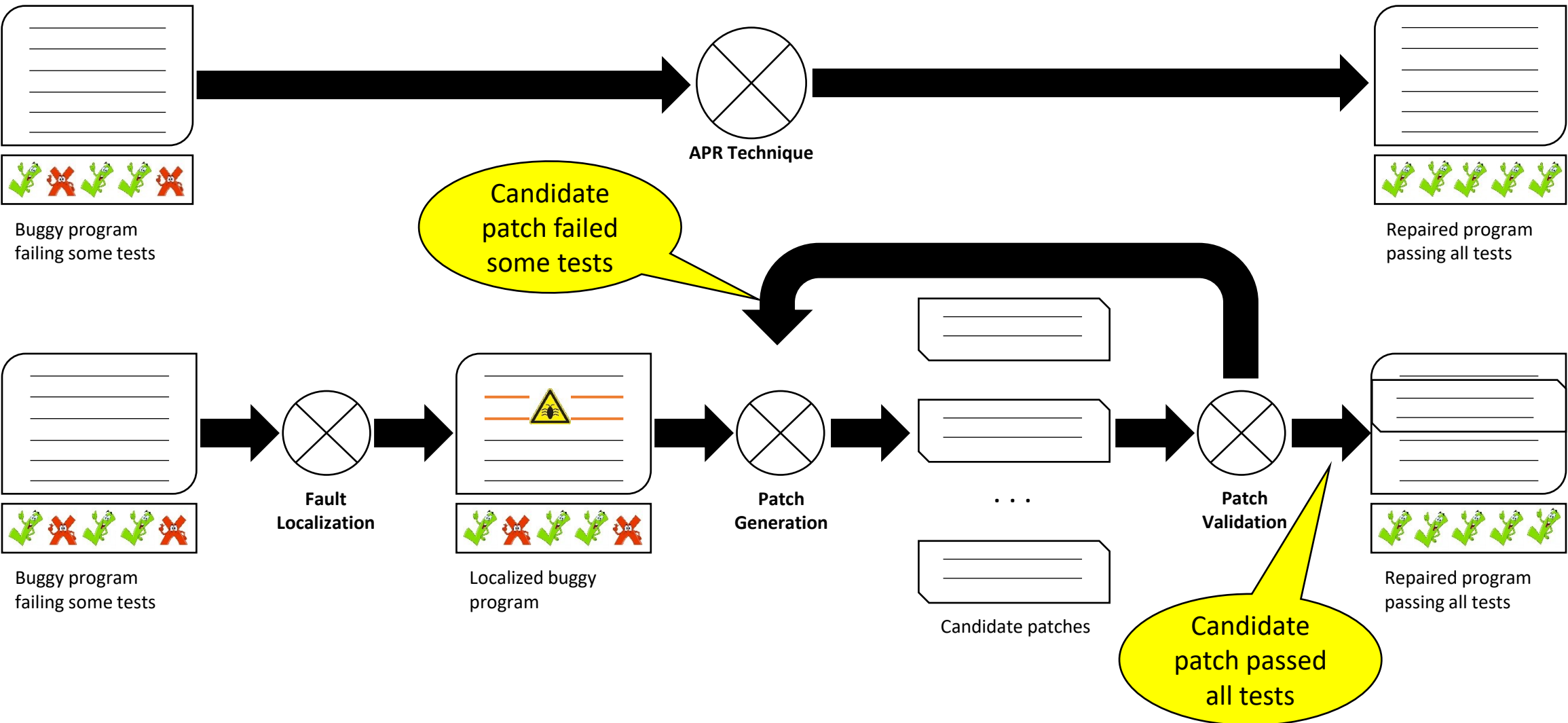## Finding and fixing software bugs automatically with SapFix and Sapienz

### Workflow (Generation)

Sapienz Auto Triage → *Bug Detected* → Trigger Patch Generator → *Triggers* → Fix Patch Generator → Validated Revision

Revert Full Diff | Revert Partial Diff | Template | Mutation

# State of the Art



http://program-repair.org

# Automatic Program Repair



APR technique

Buggy program
failing some tests

Repaired program
passing all tests
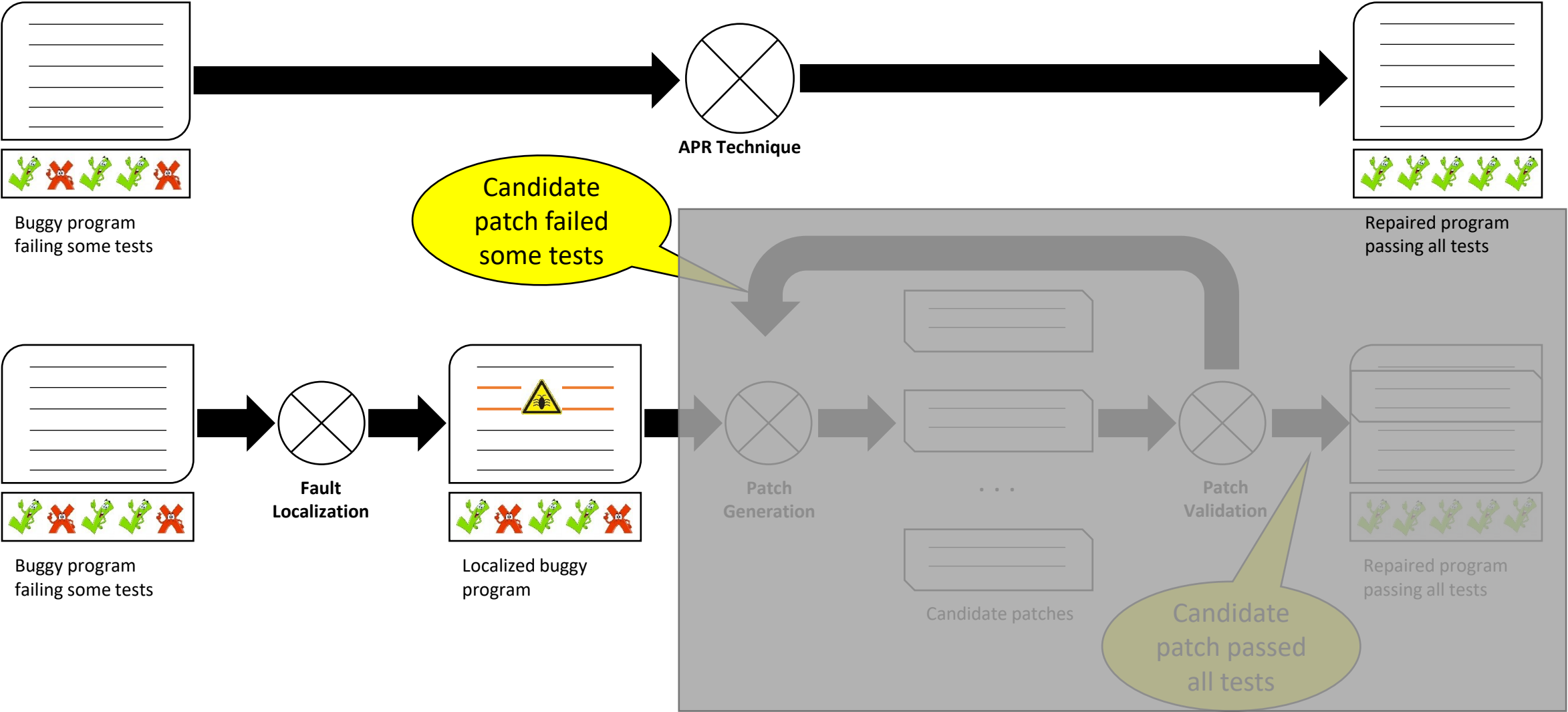
# Program Repair Process

# Program Repair Process

# Fault Localization

**Fault Localization:** *Automatically determining program elements (such as statements or methods) that are defective and cause software failure.*
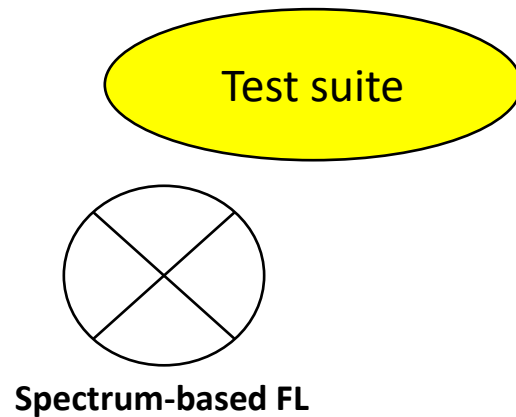
# Fault Localization Techniques

Test suite

**Spectrum-based FL**

Effect of passing and failing tests on *program element*

# Fault Localization Techniques

Test suite

Spectrum-based FL

**Effect of passing and failing tests on *program element***

**Source program**

| |
|---|
| code statement 1 |
| code statement 2 |
| ... |
| **code statement i** |
| ... |
| code statement m |

**Test suite**

| | |
|---|---|
| Test 1 | pass |
| Test 2 | fail |
| ... | ... |
| Test k | pass |
| ... | ... |
| Test n | fail |

- Count the number of passing and failing tests which execute a code statement I
- **Suspiciousness score(code statement i) = *f* (#passing tests executing stmt i, #failing tests executing stmt i)**

# Fault Localization Techniques

Test suite

Spectrum-based FL    **Mutation-based FL**

Effect of *program element* on passing and failing tests

# Fault Localization Techniques

Test suite

Spectrum-based FL    **Mutation-based FL**

Effect of *program element* on passing and failing tests

**Source program**

| code statement 1 |
| code statement 2 |
| ... |
| **code statement i** |
| ... |
| code statement m |

**Test suite**

| Test 1 | pass |
|--------|------|
| Test 2 | fail |
| ... | ... |
| Test k | pass |
| ... | ... |
| Test n | fail |

- Mutate code statement i (for e.g.,  change 'a = a + 1; ' to 'a = a - 1')
- count the total number of passing and failing tests
- **Suspiciousness score (code statement i) = $f$(total #passing tests, total #failing tests)**

# Fault Localization Techniques

Test suite

Spectrum-based FL

Mutation-based FL

Program slice

**Program Slicing FL**

Reduce program to a minimal form while still maintaining a given behavior exhibited using test suite
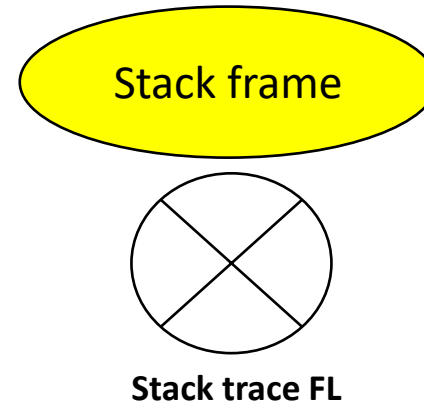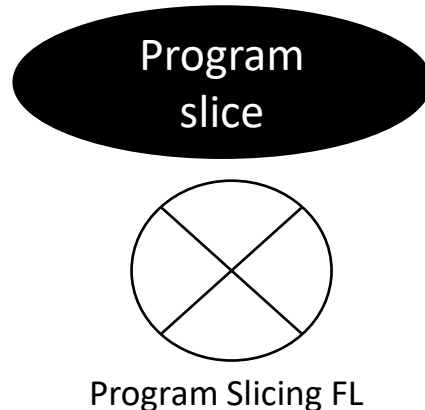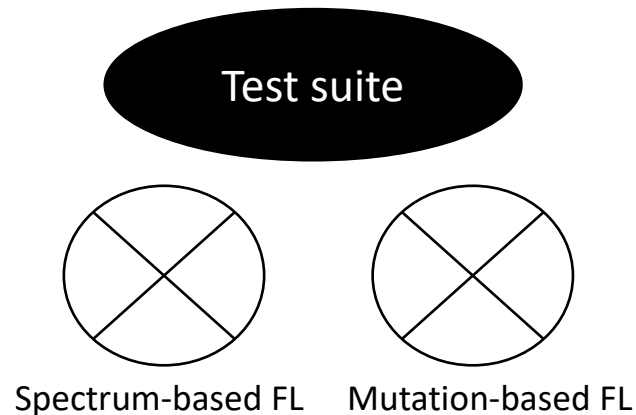
# Fault Localization Techniques

Test suite

Spectrum-based FL        Mutation-based FL

Program slice

**Program Slicing FL**

Reduce program to a minimal form while still maintaining a given behavior exhibited using test suite

**Source program**

code statement 1
...

code statement i-1
code statement i
code statement i+1

...
Code statement m

**Test suite**

| Test 1 | pass |
|--------|------|
| Test 2 | fail |
| ...    | ...  |
| Test k | pass |
| ...    | ...  |
| Test n | fail |

# Fault Localization Techniques



Test suite

Spectrum-based FL        Mutation-based FL

Program slice

Program Slicing FL

Stack frame

**Stack trace FL**

Use the list of active stack frames during execution of a program

# Fault Localization Techniques

Test suite

Spectrum-based FL     Mutation-based FL

**Source program**

If (**cond 1**){

…

} else if (**cond 2**) {

…

**Test suite**

| Test 1 | pass |
| Test 2 | fail |
| … | … |
| Test n | fail |

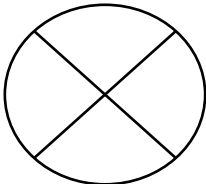Effect of predicate (conditional expression) on failing tests

- Mutate  if condition (for e.g.,  change 'if (a - b > 0) ' to ' if (a – b < 0)')
- count the total number of passing and failing tests
- **Suspiciousness score (cond i) = $f$(total #passing tests, total #failing tests)**

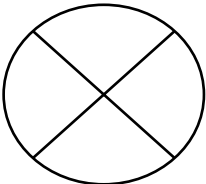Program slicing

Program Slicing FL

Stack frame

Stack trace FL

Control flow

**Predicate Switching FL**

# Fault Localization Techniques

Test Suite

Rank program elements using bug report as query

Bug report

Spectrum-based FL    Mutation-based FL

**IR-based FL**

Program Slicing

Stack Trace

Control Flow

Program Slicing FL

Stack trace FL

Predicate Switching FL

# Fault Localization Techniques

**Test Suite**

Spectrum-based FL　　Mutation-based FL

*Fault Prediction* techniques that use *code proneness* and *development history* to predict buggy program elements
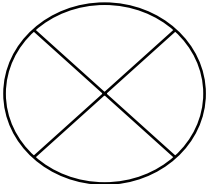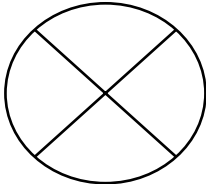
**Bug report**

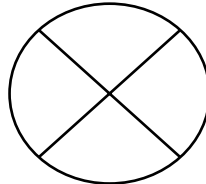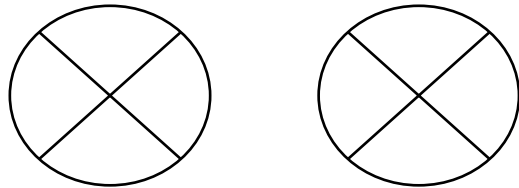IR-based FL

**Program Slicing**

Program Slicing FL

**Stack Trace**

Stack trace FL

**Control Flow**

Predicate Switching FL

# Fault Localization Techniques

Test Suite

Spectrum-based FL    Mutation-based FL

**None of these techniques is the best technique!**

**Automated program repair techniques typically use SBFL**

Bug report

IR-based FL

Program Slicing

Program Slicing FL

Stack Trace

Stack trace FL

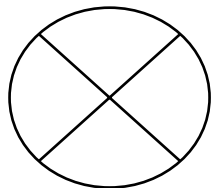Control Flow

Predicate Switching FL

# Program Repair Process



**APR Technique**

Buggy program
failing some tests

Repaired program
passing all tests

Candidate
patch failed
some tests

**Fault
Localization**

Buggy program
failing some tests

Localized buggy
program

**Patch
Generation**

. . .

Candidate patches

**Patch
Validation**

Repaired program
passing all tests

Candidate
patch passed
all tests

# Patch Generation

Three classes of repair techniques:

- **Generate-and-Validate (G&V) or Heuristics-based**
  - Use search-based software engineering (e.g., find code snippet which is similar to buggy code) to generate patch
  - E.g., GenProg, SimFix, Tbar, Arja
- **Constraint-based or Synthesis-based**
  - Builds formal constraints (e.g. SMT) and uses constraint solvers (e.g. Z3) to generate patch
  - E.g., SearchRepair, SOSRepair, Angelix
- **Learning-based**
  - Use deep learning techniques to *translate* buggy source code into patched source code (similar to translating one natural language to another)
  - E.g., Cure, Recoder, SequenceR

# Program Repair Process

# Cobra Effect

- *When an attempted solution to a problem makes the problem worse, as a type of unintended consequence.*

# Cobra Effect

- *When an attempted solution to a problem makes the problem worse, as a type of unintended consequence.*

- What has this to do with program repair?

# Patch Validation

```
1   Int triangle(int a, int b, int c) {
2       if (a <= 0 || b <= 0 || c <= 0)
3            return INVALID;
4       if (a == b && b == c)
5            return EQUILATERAL;
6       if (a == b || b != c)          // bug!
7            return ISOSCELES;
8       return SCALENE;
9   }
```

**correct patch**
```
- if (a == b || b != c)
+ if (a == b || b == c || a == c)
        return ISOSCELES
```

**overfitted/plausible patch**
```
- if (a == b || b != c)
+ if (c == 2 || c == 3)
        return ISOSCELES
```

Can we automatically identify over-fitted patches ?

| Test-id | a | b | c | Expected output | Pass/Fail |
|---------|-----|-----|-----|----------------|-----------|
| 1 | −1 | −1 | −1 | INVALID | Pass |
| 2 | 1 | 1 | 1 | EQUILATERAL | Pass |
| 3 | 2 | 2 | 3 | ISOSCELES | Pass |
| 4 | 3 | 2 | 2 | ISOSCELES | Fail |
| 5 | 2 | 3 | 2 | ISOSCELES | Fail |
| 6 | 2 | 3 | 4 | SCALENE | Fail |

Example source: Le Goues et al., Automated program repair. CACM, Nov 2019

# How to evaluate patch "correctness"?
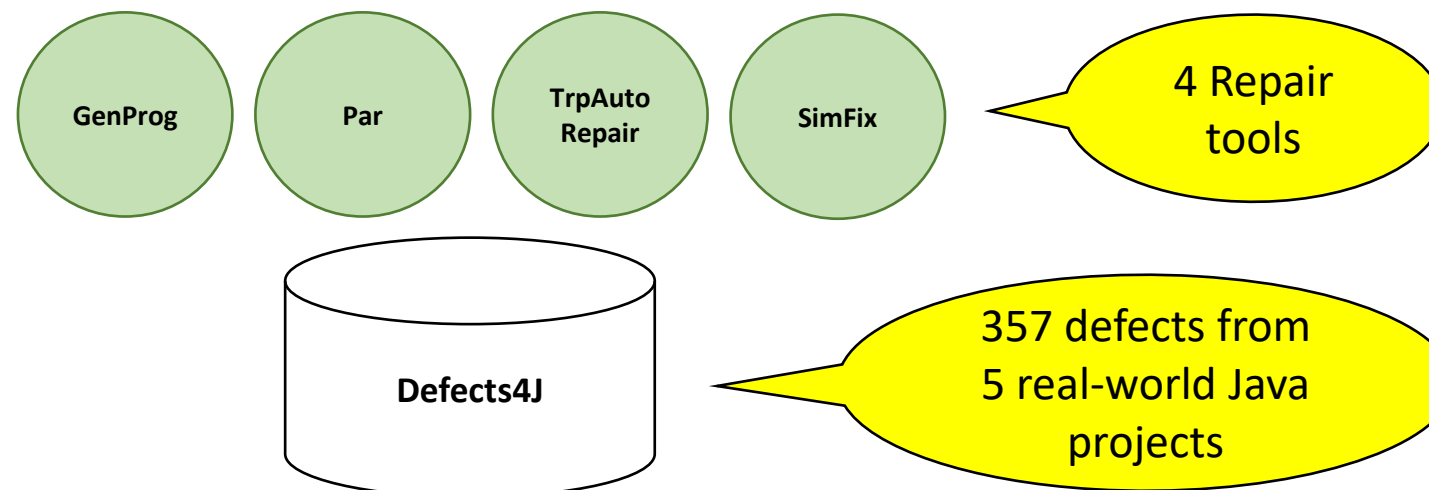
- Use automatically generated *high-quality* evaluation test-suite to measure the patch quality/correctness

- Patch Quality = #tests passed/total #tests

# Quality of Heuristics Repair Techniques



Manish Motwani, Mauricio Soto, Yuriy Brun, René Just, and Claire Le Goues, Quality of Automated Program Repair on Real-World Defects, IEEE TSE, 2022.

# Quality of Heuristics Repair Techniques

**RQ1**: Do G&V techniques produce patches for real-world Java defects?

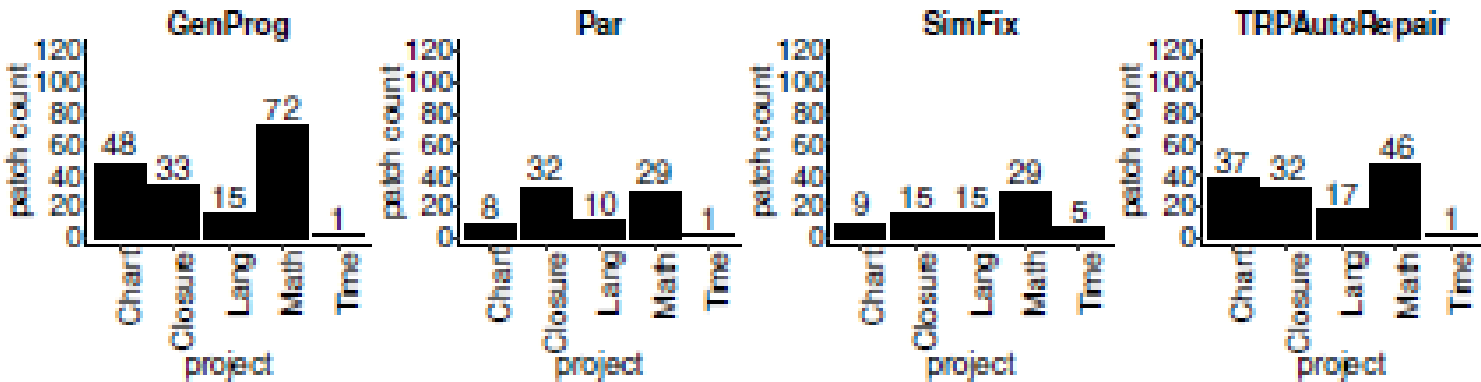Manish Motwani, Mauricio Soto, Yuriy Brun, René Just, and Claire Le Goues, Quality of Automated Program Repair on Real-World Defects, IEEE TSE, 2022.

# Quality of Heuristics Repair Techniques

**RQ1**: Do G&V techniques produce patches for real-world Java defects?

| technique | patches | | defects |
|---|---|---|---|
| | total | unique | patched |
| GenProg | 585 (8.2%) | 255 | 49 (13.7%) |
| Par | 288 (4.0%) | 107 | 38 (10.6%) |
| SimFix | 76 (21.3%) | 73 | 68 (19.0%) |
| TRPAutoRepair | 513 (7.2%) | 199 | 44 (12.3%) |
| total | 1,462 ( 6.7%) | 634 | 106 (29.7%) |



**RA1**: Yes, although less often than for C defects (19% vs ~50%).

Manish Motwani, Mauricio Soto, Yuriy Brun, René Just, and Claire Le Goues, Quality of Automated Program Repair on Real-World Defects, IEEE TSE, 2022.
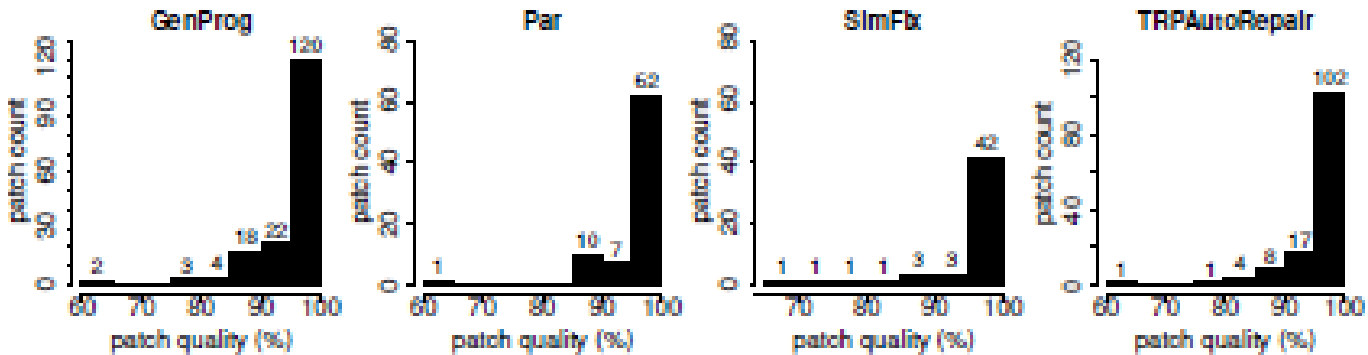
# Quality of Heuristics Repair Techniques

**RQ2**: How often and how much do the patches produced by repair techniques overfit to the developer-written test suite and fail to generalize to the evaluation test suite?

Manish Motwani, Mauricio Soto, Yuriy Brun, René Just, and Claire Le Goues, Quality of Automated Program Repair on Real-World Defects, IEEE TSE, 2022.
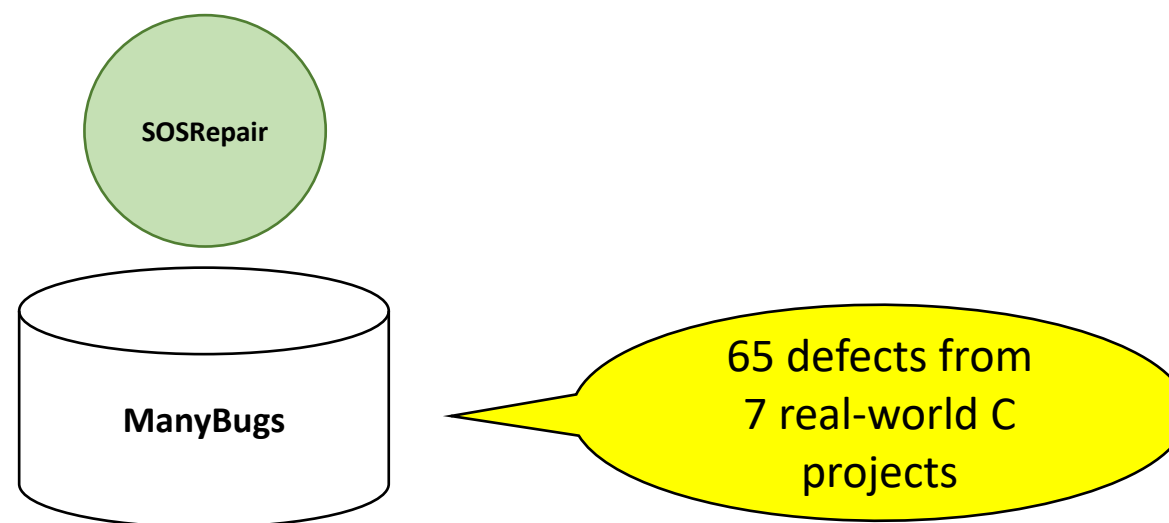
# Quality of Heuristics Repair Techniques

**RQ2**: How often and how much do the patches produced by repair techniques overfit to the developer-written test suite and fail to generalize to the evaluation test suite?

| technique | minimum | patch quality | | maximum | 100%-quality patches |
|---|---|---|---|---|---|
| | | mean | median | | |
| GenProg | 64.8% | 95.7% | 98.4% | 100.0% | 24.3% |
| Par | 64.8% | 96.1% | 98.5% | 100.0% | 13.8% |
| SimFix | 65.0% | 96.3% | 99.9% | 100.0% | 46.1% |
| TrpAutoRepair | 64.8% | 96.4% | 98.4% | 100.0% | 19.5% |



**RA2**: Often. Only between 13.8% and 46.1% of the patches pass 100% of evaluation test suite.

Manish Motwani, Mauricio Soto, Yuriy Brun, René Just, and Claire Le Goues, Quality of Automated Program Repair on Real-World Defects, IEEE TSE, 2022.

# Quality of Synthesis-based Repair Technique



SOSRepair

ManyBugs

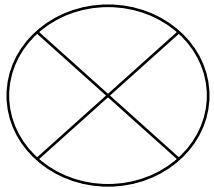65 defects from 7 real-world C projects

*SOSRepair patches 22 (34%) out of 65 defects and of these 22 patches, 9 (41%) pass all independent tests*

*Key-Finding: Manually improving fault localization allows SOSRepair to patch 23 (35%) of the defects, of which 16 (70%) pass all independent tests*
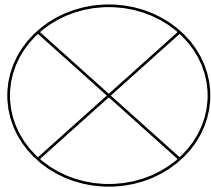
Afsoon Afzal, Manish Motwani, Kathryn T. Stolee, Yuriy Brun, and Claire Le Goues, SOSRepair: Expressive Semantic Search for Real-World Program Repair, IEEE TSE, 2021

# How can we improve the quality of repair techniques?

# Idea-1: Improve Fault Localization

Test Suite

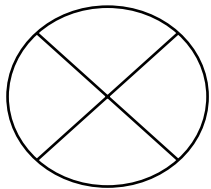Spectrum-based FL          Mutation-based FL

Combine multiple techniques which use different software artifacts to identify buggy program elements.
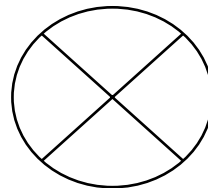
Bug report
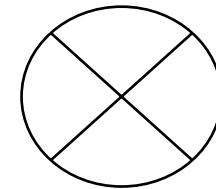
IR-based FL

Program Slicing

Program Slicing FL

Stack Trace

Stack trace FL

Control Flow

Predicate Switching FL

# Idea-2: Improve Patch Generation

Target specific class of defects for e.g., fixing bugs in if conditionals or fixing null pointer exceptions

Use learned fix patterns instead of hard-coded ones

Use novel search strategies to produce the correct patch first

Use novel test sampling strategies to test candidate patches on a sample of test suite and reduce computation time and resources

Use information from other artifacts such as contracts/specifications and bug reports

# Idea-3: Improve Patch Validation

Can we improve tests using information from artifacts that are available in real scenario when defect occurs?

Yes! we can use Natural Language Processing (NLP) techniques to generate executable tests from natural-language software artifacts
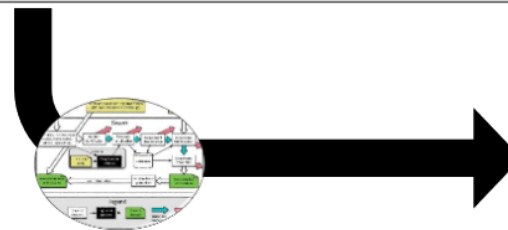


swami.cs.umass.edu

# What did we learn?

- Three step process of automatic program repair
  - Fault localization
  - Patch generation
  - Patch validation
- Patch overfitting (repair quality) is a real concern
  - More than 50% of the patches produced by repair techniques overfit.
- Ideas to improve the repair quality by:
  - Improving fault localization
  - Improving patch generation
  - Improving patch validation
- For the latest updated research in program repair, visit http://program-repair.org/
- You will learn how GenProg repairs real-world bugs in HW2.

# Announcements

- HW2 is released today and is due on Wednesday, May 29, 11:59 PM
- Next class, we will have project plan presentations.
  - Each project group will give **20 min** presentation.
  - The presentation should **cover all the aspects of your project except results**.
  - Remember to **use key principles** while organizing your talk and text in your report.
  - Utilize office hours tomorrow from 2-3 PM if you want to discuss something.