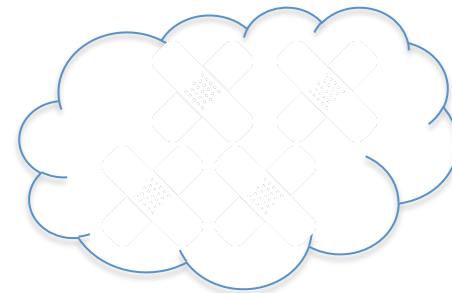# An Analysis of Patch Plausibility and Correctness of Generate-And-Validate Patch Generation System

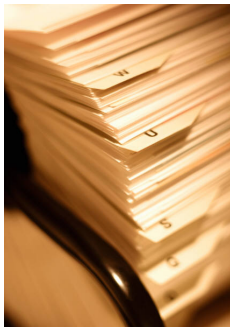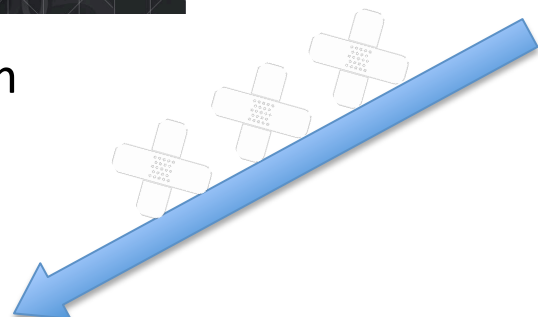**Zichao Qi**, Fan Long, Sara Achour, and Martin Rinard

MIT CSAIL

# Generate-And-Validate Patch Generation Systems

Buggy Program

Candidate patch space

Test suite of test cases

# Generate-And-Validate Patch Generation Systems

- **GenProg** – Genetic Programming
  1. C. L. Goues, M. Dewey-Vogt, S. Forrest, and W. Weimer. *A systematic study of automated program repair: Fixing 55 out of 105 bugs for $8 each.* **ICSE 2012**
  2. W. Weimer, T. Nguyen, C. Le Goues, and S. Forrest. *Automatically finding patches using genetic programming.* **ICSE 2009**
  3. S. Forrest, T. Nguyen, W. Weimer, and C. Le Goues. *A genetic programming approach to automated software repair.* **GECCO 2009**
  4. C. Le Goues, T. Nguyen, S. Forrest, and W. Weimer. Genprog: *A generic method for automatic software repair.* **Software Engineering, IEEE Transactions on 38(1), 2012**

- **AE** – Adaptive Search
  1. W. Weimer, Z. P. Fry, and S. Forrest. *Leveraging program equivalence for adaptive program repair: Models and first results.* **ASE 2013**

- **RSRepair** – Random Search
  1. Y. Qi, X. Mao, Y. Lei, Z. Dai, and C. Wang. *The strength of random search on automated program repair.* **ICSE 2014**

# All of them report impressive results

|  | GenProg | AE | RSRepair |
|---|---|---|---|
| **Benchmark Defects** | 105 | 105 | 24 |
| **Reported Fixed Defects** | 55 | 54 | 24 |

- Patches generated by these systems are different from human written patch
- No systematic analysis

# We analyze the reported patches for these systems

**Plausible?**
Produce correct outputs for all
test cases in the test suite

*All generated patches should
be plausible*

# Plausibility

| | GenProg | AE | RSRepair |
|---|---|---|---|
| **Benchmark Defects** | 105 | 105 | 24 |
| **Reported Fixed Defects** | 55 | 54 | 24 |
| **Defects With Plausible Patches** | 18 | 27 | 10 |

- Reason - Weak Proxy
  - Patch evaluation does not check for correct output
  - php, libtiff – check exit code, *not output*
  - Accepted php patch: main(){ exit(0); }

# Analysis of the reported patches for these systems

**Plausible?**
Produce correct outputs for all test cases in the test suite

**Majority of the patches are not plausible**

**Correct?**
Eliminate the defect

*Passing test suite != correctness*

# Correctness

| | GenProg | AE | RSRepair |
|---|---|---|---|
| **Benchmark Defects** | 105 | 105 | 24 |
| **Reported Fixed Defects** | 55 | 54 | 24 |
| **Defects With Plausible Patches** | 18 | 27 | 10 |
| **Defects With Correct Patches** | 2 | 3 | 2 |

Developed new test cases that expose defects for all plausible but incorrect patches

# GenProg Statistics



- 2 Correct
- 16 Plausible but Incorrect
- 37 Implausible

**Stronger Test Suites?**
Will GenProg generate correct patches given new test cases that eliminate incorrect patches?

**Fixed Test Scripts?**
Will GenProg generate plausible patches given fixed patch evaluation scripts?

# Analysis of the reported patches for these systems

**Plausible?**
Produce correct outputs for all test cases in the test suite

**Majority of the patches are not plausible**

**Correct?**
Eliminate the defect

**The overwhelming majority of the patches are not correct**

**Do stronger test suites help?**

**Rerun GenProg with fixed patch evaluation scripts and new test cases that eliminate incorrect patches**

# Reexecution of GenProg on Remaining 103 Defects

**First Reexecution**
Fixed patch evaluation
New test cases

Patches for 2 defects

**Second Reexecution**
2 additional test cases

Patches for 0 defects

# Why?

- Developer patches are not in GenProg search space
- GenProg search space may not contain *any* correct patch for these 103 defects
- May need richer search space to generate correct patches

# Bottom Line For GenProg

- Rerun GenProg with
  - Fixed test scripts
  - Stronger test suites
- GenProg generates patches for only 2 of 105 defects (both patches are correct)

# Examples of Correct GenProg Patch(1/2)

## Developer

```
1    -if (y < 1000) {
2    -       PyObject *accept = PyDict_GetItemString(moddict,
3    -                                               "accept2dy
4    -       if (accept != NULL) {
5    -           int acceptval = PyObject_IsTrue(accept);
6    -           if (acceptval == -1)
7    -               return 0;
8    -           if (acceptval) {
9    -               if (0 <= y && y < 69)
10   -                   y += 2000;
11   -               else if (69 <= y && y < 100)
12   -                   y += 1900;
13   -               else {
14   -                   PyErr_SetString(PyExc_ValueError,
15   -                                   "year out of range");
16   -                   return 0;
17   -               }
18   -               if (PyErr_WarnEx(PyExc_DeprecationWarning
19   -                   "Century info guessed for a 2-digit ye
20   -                   return 0;
21   -           }
22   -       }
23   -       else
24   -           return 0;
25   -}
26   p->tm_year = y - 1900;
27   p->tm_mon--;
28   p->tm_wday = (p->tm_wday + 1) % 7;
```

## GenProg

```
1    - if (y < 1000) {
2    -       tmp___0 = PyDict_GetItemString(moddict, "accept2dyear");
3    -       accept = tmp___0;
4    -       if ((unsigned int )accept != (unsigned int )((void *)0)) {
5    -           tmp___1 = PyObject_IsTrue(accept);
6    -           acceptval = tmp___1;
7    -           if (acceptval == -1) {
8    -               return (0);
9    -           } else {
10
11   -           }
12
13   -           if (acceptval) {
14   -               if (0 <= y) {
15   -                   if (y < 69) {
16   -                       y += 2000;
17   -                   } else {
18   -                       goto _L;
19   -                   }
20   -               } else {
21   -                   _L: /* CIL Label */
22   -                   if (69 <= y) {
23   -                       if (y < 100) {
24   -                           y += 1900;
25   -                       } else {
26   -                           PyErr_SetString(PyExc_ValueError,
27   -                                           "year out of range");
28   -                           return (0);
29   -                       }
30   -                   } else {
31   -                       PyErr_SetString(PyExc_ValueError,
32   -                                       "year out of range");
33   -                       return (0);
34   -                   }
35   -               }
36
37   -               tmp___2 = PyErr_WarnEx(PyExc_DeprecationWarning,
38   -                           "Century info guessed for a 2-digit year.", 1);
39   -               if (tmp___2 != 0) {
40   -                   return (0);
41   -               } else {
42
43   -               }
44   -           } else {
45
46   -           }
47   -       } else {
48   -           return (0);
49   -       }
50   - } else {
51
52   - }
53   p->tm_year = y - 1900;
54   (p->tm_mon) --;
55   p->tm_wday = (p->tm_wday + 1) % 7;
```

# Examples of Correct GenProg Patch(2/2)

Developer

GenProg

# All Correct Patches Simply Delete Code

# Semantic Analysis

- Analyze all the plausible patches
- Determine if patch is equivalent to <span style="color:red">single functionality deletion</span> modification
- Results
  - GenProg: 14/18
  - AE: 22/27
  - RSRepair: 8/10

# We found a common scenario

- A negative test case exposes the defect
  - Feature is otherwise unexercised
  - The patch simply deletes the functionality
    - Introduces new security vulnerabilities
      (buffer overflows)
    - Disables critical functionality
      (gzip cannot decompress non-zero files)

- Weak test suites
  - May be appropriate for human developers
  - *May not be* appropriate for automatic patch generation systems (at least not by themselves)

If all these patches simply delete functionality

Why not build a patch generation system that ONLY deletes functionality?

# We present Kali

- Automatic patch generation system

- Consider the search space that consists of only patches that remove functionality

# Experimental Results of Kali

| | GenProg | AE | RSRepair | Kali |
|---|---|---|---|---|
| **Benchmark Defects** | 105 | 105 | 24 | 105 |
| **Reported Fixed Defects** | 55 | 54 | 24 | |
| **Defects With Plausible Patches** | 18 | 27 | 10 | 27 |
| **Defects With Correct Patches** | 2 | 3 | 2 | 3 |

- Kali is as good as previous systems
  - Much simpler
  - Not need to know the source code file to repair
- Can pinpoint the defective code
- Can provide insight into important defect characteristics.

# Experimental Results of Kali

| | GenProg | AE | RSRepair | Kali |
|---|---|---|---|---|
| **Benchmark Defects** | 105 | 105 | 24 | 105 |
| **Reported Fixed Defects** | 55 | 54 | 24 | |
| **Defects With Plausible Patches** | 18 | 27 | 10 | 27 |
| **Defects With Correct Patches** | 2 | 3 | 2 | 3 |

# Is Automatic Patch Generation A
## Total Failure?

# NO!

# Path To Success

- Richer search spaces
- More efficient search algorithms
- Incorporate additional sources of information
  - Correct code from other applications
  - Learned characteristics of human patches
  - Learned invariants
  - Specifications

# Promising directions

- Learn invariant from correct execution
  - ClearView: *J. H. Perkins, S. Kim, S. Larsen, S. Amarasinghe, J. Bachrach, M. Carbin, C. Pacheco, F. Sherwood, S. Sidiroglou, G. Sullivan, et al.* *Automatically patching errors in deployed software.* *SOSP 2009.*
  - Patches security vulnerabilities in 9 of 10 defects
  - At least 4 patches are correct
- Solvers
  - NOPOL: F. *DeMarco, J. Xuan, D. Le Berre, and M. Monperrus.* *Automatic repair of buggy if conditions and missing preconditions with smt.* *CSTVA 2014*
  - SemFix: *H. D. T. Nguyen, D. Qi, A. Roychoudhury, and S. Chandra. Semfix:* *Program repair via semantic analysis.* *ICSE 2013*

# Promising directions

- ## Specifications
  - Autofix-E: *Yu Pei, Carlo A. Furia, Martin Nordio, Yi Wei, Andreas Zeller, and Bertrand Meyer. Automated Fixing of Programs with Contracts. IEEE Transactions on Software Engineering, 2014.*
  - *Etienne Kneuss, Manos Koukoutos and Viktor Kuncak. Deductive Program Repair. CAV 2015*

- ## Correctness evaluation
  - *Thomas Durieux, Matias Martinez, Martin Monperrus, Romain Sommerard, Jifeng Xuan. Automatic Repair of Real Bugs: An Experience Report on the Defects4J Dataset. Technical report 1505.07002, Arxiv, 2015*

- ## Code from another application
  - CodePhage: *S. Sidiroglou, E. Lahtinen, F. Long, and M. Rinard. Automatic error elimination by multi-application code transfer. PLDI 2015*

# Promising Directions

| | **SPR**<br>Staged condition synthesis | **Prophet**<br>Learn from successful patches |
|---|:---:|:---:|
| **Benchmark Defects** | **105** | **105** |
| **Defects With Plausible Patches** | **41** | **42** |
| **Defects With Correct Patches** | **11** | **15** |

SPR: *F. Long and M. Rinard. Staged program repair in SPR. To appear in ESEC-FSE 2015*

Prophet: *F. Long and M. Rinard. Prophet: Automatic patch generation via learning from successful human patches. Under submission*

# Take Aways

- Facts about GenProg/AE/RSRepair
  - These systems fix 2/3/2 of 105 bugs (not 55/54/24)
  - Errors in test scripts and weak test suites
  - Fixed test scripts and stronger test suites do not help
- Paths to success
  - Richer search spaces
  - More efficient search algorithms
  - Incorporate additional sources of information
    - Correct code from other applications (CodePhage)
    - Learned characteristics of human patches (Prophet)
    - Learned invariants (ClearView)
    - Specifications (AutoFixE, Deductive Repair)

# Summary

- Evaluation of GenProg, AE and RSRepair
  - Incorrect results
  - Equivalent to functionality elimination
  - Stronger test suites do not help
- Kali
  - Functionality elimination system
  - Help developer better understand the bug

# Path to Success for the automatic patch generation systems

- Richer search spaces
- More efficient search algorithms
- Incorporate additional sources of information
  - Correct code from other applications
  - Learned characteristics of human patches
  - Learned invariants
- Better patch evaluation
  - Correctness
  - Understand the negative effects

Questions?