## Upcoming

- Homework 3 due April 18
- Literature review due today April 11

1

## Repairing Automated Repair

2

## Repairing Automated Repair

3

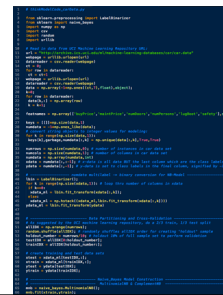## Cobra Effect

4

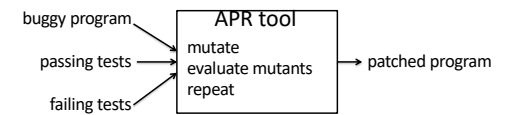## What do cobras have to do with automated program repair?

repairing python programs?

5

## Automated Program Repair

basic idea:

buggy program → | APR tool
passing tests → | mutate
failing tests → | evaluate mutants
repeat | → patched program

6

## Slide 7

### the many repair tools

ClearView [Perkinds et al. 2009] GenProg [Weimer et al. 2009]

Prophet [Long and Rinard 2015]   SPR [Long and Rinard 2015]

TDS [Perelman et al. 2014]

Par [Kim et al. 2013]   AE [Weimer et al. 2013]

SemFix [Nguyen et al. 2013]   AutoFix-E [Wei et al. 2010]

[Carzaniga et al. 2010]   [Carzaniga et al. 2013]

[Jin et al. 2011]   Coker and Hafiz et al. 2013]

[Debroy and Wong et al. 2010]   [Lin and Ernst et al. 2004]

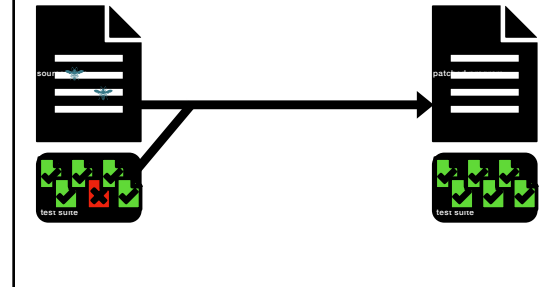[Forrest et al. 2009]   [Novark et al. 2007]   [Demsky et al. 2006]

7
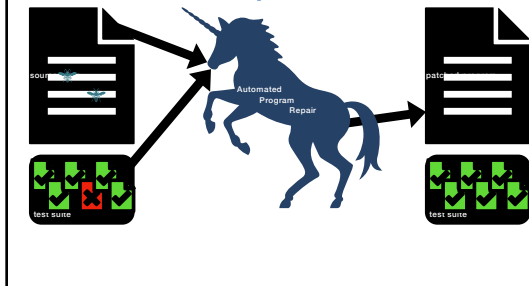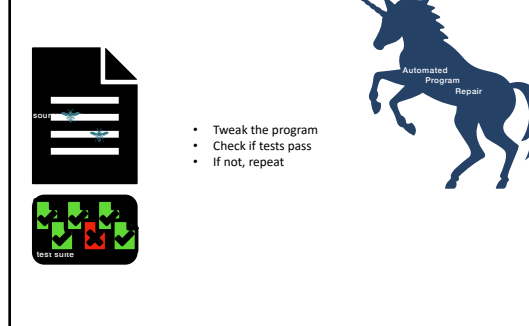
## Slide 8

### Program repair techniques



8

## Slide 9

# Automated program repair



9

## Slide 10

# Automated program repair



10

## Slide 11

### Program repair techniques



- Tweak the program
- Check if tests pass
- If not, repeat

11

## Slide 12

### Potential problem

buggy program → **APR tool** (mutate, evaluate mutants, repeat) → patched program

passing tests →

failing tests →

the patched program may pass all given tests, but break other functionality

12

# COMPUTE THE MEDIAN OF THREE NUMBERS

13

```
int median(int a, int b, int c) {
  int result;
  if ((b<=a && a<=c) ||
      (c<=a && a<=b))
    result = a;
  if ((a<b && b <= c) ||
      (c<=b && b<a))
    result = b;
  if ((a<c && c<b) ||
      (b<c && c<a))
    result = c;
  return result;
}
```

14

```
int median(int a, int b, int c) {
  int result = 0;
  if ((b<=a && a<=c) ||
      (c<=a && a<=b))
    result = a;
  if ((a<b && b <= c) ||
      (c<=b && b<a))
    result = b;
  if ((a<c && c<b) ||
      (b<c && c<a))
    result = c;
  return result;
}
```

15

```
int median(int a, int b, int c) {
  int result = 0;
  if ((b<=a && a<=c) ||
      (c<=a && a<=b))
    result = a;
  if ((a<b && b <= c) ||
      (c<=b && b<a))
    result = b;
  if ((a<c && c<b) ||
      (b<c && c<a))
    result = c;
  return result;
}
```

16

```
int median(int a, int b, int c) {
  int result = 0;
  if ((b<=a && a<=c) ||
      (c<=a && a<=b))
    result = a;
  if ((a<b && b <= c) ||
      (c<=b && b<a))
    result = b;
  if ((a<c && c<b) ||
      (b<c && c<a))
    result = c;
  return result;
}
```

17

```
int median(int a, int b, int c) {
  int result = 0;
  if ((b<=a && a<=c) ||
      (c<=a && a<=b))
    result = a;
  if ((a<b && b <= c) ||
      (c<=b && b<a))
    result = b;
  if ((a<c && c<b) ||
      (b<c && c<a))
    result = c;
  return result;
}
```

18

```
int median(int a, int b, int c) {
  int result = 0;
  if ((b<=a && a<=c) ||
      (c<=a && a<=b))
    result = a;
  if ((a<b && b <= c) ||
      (c<=b && b<a))
    result = b;
  if ((a<c && c<b) ||
      (b<c && c<a))
    result = c;
  return result;
}
```

19

```
int median(int a, int b, int c) {
  int result = 0;
  if ((b<=a && a<=c) ||
      (c<=a && a<=b)
    result = a;
  if ((a<b && b <= c) ||
      (c<=b && b<a))
    result = b;
  if ((a<c && c<b) ||
      (b<c && c<a))
    result = c;
  return result;
}
```

20

```
int median(int a, int b, int c) {
  int result = 0;
    ((b<=a && a<=c) ||
     (c<=a && a<=b))
    result = a;
    ((a<b && b <= c) ||
     (c<=b && b<a))
    result = b;
    ((a<c && c<b) ||
     (b<c && c<a))
    result = c;
  return result;
}
```

21

```
int med_broken(int a, int b, int c) {
  int result;
  if ((a==b) || (a==c) ||
      (b<a && a<c) ||
      (c<a && a<b))
    result = a;
  else if ((b==c) || (a<b && b<c) ||
           (c<b && b<a))
    result = b;
  else if (a<c && c<b)
    result = c;
  return result;
}
```

22

```
int med_broken(int a, int b, int c) {
  int result;
  if ((a==b) || (a==c) ||
      (b<a && a<c) ||
      (c<a && a<b))
    result = a;
  else if ((b==c) || (a<b && b<c) ||
           (c<b && b<a))
    result = b;
  else if (a<c && c<b)
    result = c;
  return result;
}
```

23

```
int med_broken(int a, int b, int c) {
  int result;
  if ((a==b) || (a==c) ||
      (b<a && a<c) ||
      (c<a && a<b))
    result = a;
  else if ((b==c) || (a<b && b<c) ||
           (c<b && b<a))
    result = b;
  else if (a<c && c<b)
    result = c;
  return result;
}
```

24

4

**Slide 25**

```
int med_broken(int a, int b, int c) {
  int result;
  if ((a==b) || (a==c) ||
       (b<a && a<c) ||
       (c<a && a<b))
    result = a;
  else if ((b==c) || (a<b && b<c) ||
            (c<b && b<a))
    result = b;
  else if (a<c && c<b)
    result = c;
  return result;
}
```

| Input | Expected | Pass? |
|-------|----------|-------|
| 0,0,0 | 0 | ✓ |
| 2,0,1 | 1 | ✗ |
| 0,0,1 | 0 | ✓ |
| 0,1,0 | 0 | ✓ |
| 0,2,1 | 1 | ✓ |
| 0,2,3 | 2 | ✓ |

25

**Slide 26**

```
int med_broken(int a, int b, int c) {
  int result;
  if ((a==b) || (a==c) ||
       (b<a && a<c) ||
       (c<a && a<b))
    result = a;
  if (b < a)
    result = c;
  else if (b<a) (b==c) || (a<b && b<c) ||
            (c<b && b<a))
    result = b;
  else if (a<c && c<b)
    result = c;
  return result;
}
```

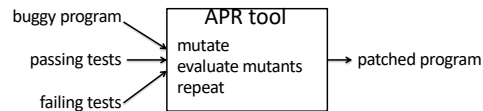| Input | Expected | Pass? |
|-------|----------|-------|
| 0,0,0 | 0 | ✓ |
| 2,0,1 | 1 | ✗ |
| 0,0,1 | 0 | ✓ |
| 0,1,0 | 0 | ✓ |
| 0,2,1 | 1 | ✓ |
| 0,2,3 | 2 | ✓ |

26

**Slide 27**

```
int med_broken(int a, int b, int c) {
  int result;
  if ((a==b) || (a==c) ||
       (b<a && a<c) ||
       (c<a && a<b))
    result = a;
  if (b < a)
    result = c;
  if (b<a) (b==c) || (a<b && b<c) ||
            (c<b && b<a))
    result = b;
  if (a<c && c<b)
    result = c;
  return result;
}
```

| Input | Expected | Pass? |
|-------|----------|-------|
| 0,0,0 | 0 | ✓ |
| 2,0,1 | 1 | ✗ |
| 0,0,1 | 0 | ✓ |
| 0,1,0 | 0 | ✓ |
| 0,2,1 | 1 | ✓ |
| 0,2,3 | 2 | ✓ |

27

**Slide 28**

```
int med_broken(int a, int b, int c) {
  int result;
  if ((a==b) || (a==c) ||
       (b<a && a<c) ||
       (c<a && a<b))
    result = a;
  if (b < a)
    result = c;
  else if (b<a) (b==c) || (a<b && b<c) ||
            (c<b && b<a))
    result = b;
  else if (a<c && c<b)
    result = c;
  return result;
}
```

| Input | Expected | Pass? |
|-------|----------|-------|
| 0,0,0 | 0 | ✓ |
| 2,0,1 | 1 | ✓ |
| 0,0,1 | 0 | ✓ |
| 0,1,0 | 0 | ✓ |
| 0,2,1 | 1 | ✓ |
| 0,2,3 | 2 | ✓ |

28

**Slide 29**

```
int med_broken(int a, int b, int c) {
  int result;
  if ((a==b) || (a==c) ||
       (b<a && a<c) ||
       (c<a && a<b))
    result = a;
  if ((b==c) || (a<b && b<c) ||
            (c<b && b<a))
    result = b;
  if (a<c && c<b)
    result = c;
  return result;
}
```

| Input | Expected | Pass? |
|-------|----------|-------|
| 2,6,8 | 6 | ✓ |
| 2,8,6 | 6 | ✓ |
| 6,2,8 | 6 | ✓ |
| 6,8,2 | 6 | ✓ |
| 8,2,6 | 6 | ✗ |
| 8,6,2 | 6 | ✓ |
| 9,9,9 | 9 | ✓ |

29

**Slide 30**

```
int med_broken(int a, int b, int c) {
  int result;
  if ((a==b) || (a==c) ||
       (b<a && a<c) ||
       (c<a && a<b))
    result = a;
  if (b < a)
    result = c;
  else if (b<a) (b==c) || (a<b && b<c) ||
            (c<b && b<a))
    result = b;
  else if (a<c && c<b)
    result = c;
  return result;
}
```

| Input | Expected | Pass? |
|-------|----------|-------|
| 0,0,0 | 0 | ✓ |
| 2,0,1 | 1 | ✓ |
| 0,0,1 | 0 | ✓ |
| 0,1,0 | 0 | ✓ |
| 0,2,1 | 1 | ✓ |
| 0,2,3 | 2 | ✓ |
| 2,6,8 | 6 | ✓ |
| 2,8,6 | 6 | ✓ |
| 6,2,8 | 6 | ✗ |
| 6,8,2 | 6 | ✓ |
| 8,2,6 | 6 | ✓ |
| 8,6,2 | 6 | ✗ |
| 9,9,9 | 9 | ✓ |

30

## Potential solution



buggy program → | APR tool<br>mutate<br>evaluate mutants<br>repeat | → patched program

passing tests →

failing tests →

Use an independent test suite to measure quality of the patch

31

---

## Focus of prior evaluations

- Most evaluations are interested in whether tools work
  - produce patches
- Some interest in other factors
  - human acceptance of patches [Durieux et al. 2015] [Fry et al. 2012] [Kim et al. 2013]
  - plausibility [Qi et al. 2015]
  - …but these don't fully assess functional correctness
- No evaluations test functional correctness of repair outputs independently of repair inputs

32

---

## What do we need?

- We need bugs with 2 test suites
  - and the test suites need to be good

### Why?

- it's hard enough to find one good test suite, good luck finding programs with two

33

---

## Make your own!

http://repairbenchmarks.cs.umass.edu

998 student-written buggy C programs

- simple (very small)
- have 2 test suites
  - white-box (generated by KLEE)
  - black-box (written by instructor)

Some programs fail some wb tests, others bb tests, others, some of both

34

---

## RQ1:
## What is the base incidence of overfitting?

Give a repair tool the buggy program and the black-box test suite, try to repair it, see what fraction of the white-box tests the patches pass.

35

---

## RQ1:
## What is the base incidence of overfitting?

but first, how often can we actually generate patches?

| repair tool | patch production % |
| --- | --- |
| GenProg | 466/778 = 59.9% |
| TrpAutoRepair | 444/778 = 57.1% |

36

## RQ1:
## What is the base incidence of overfitting?



37

## RQ2: What effect do pre-repair test failures have on overfitting?



**Programs that fail more tests before repair still fail more tests after repair**

38

## RQ2: What effect do pre-repair test failures have on overfitting?



**Repair is at best unlikely to improve correctness, at worst likely to worsen it**

39

## RQ3: What effect does test suite coverage have on overfitting?

- Randomly sample 25%, 50%, and 75% of passing and failing tests for each buggy program
- Attempt to repair programs
  – with each level of test coverage
- If a repair is found, measure correctness of repair

40

## RQ3: What effect does test suite coverage have on overfitting?



**Lower test suite coverage leads to more overfitting**

41

## RQ4: What effect does test suite provenance have on overfitting?

- So far, all experiments have used human-written *black-box* tests to build repairs
- Switch to using KLEE-generated *white-box* tests
- Attempt to repair programs
- If a repair is found, measure correctness of repair
  – this time with *black-box* tests

42

7

## Slide 43

### RQ4: What effect does test suite provenance have on overfitting?



**Automatically generated tests produced significantly buggier repairs compared to human-written tests**

43

## Slide 44

### RQ4: Do tools do better than novices?



44

## Slide 45

### Summary of that study

- Overfitting is a real concern
  - median patch for either tool passed only 75% of evaluation suite
- Overfitting is hard to avoid
  - minimization doesn't help on this dataset
  - N-version voting only works in extreme cases
- Program repair is harder for buggier programs, but likely to break more correct programs
- Novice developers don't significantly beat repair tools

45

## Slide 46

### How well does APR work?



46

## Slide 47

### How well does APR work?

Results on 11 tools, 2,141 bugs:
- Generate patches for 15-213 bugs
- Evaluations overfit to their benchmarks

47

## Slide 48

### How well does APR work?

11 tools generated patches on 15-213 out of 2,141 bugs.

48

### Slide 49

## How well does APR work?

- Evaluated 4 techniques
  - GenProg
  - Par
  - TrpAutoRepair
  - SimFix
- Measured patch quality
- Measured what affects patch quality

*Quality of Automated Program Repair on Real-World Defects*

49

### Slide 50

## Quality vs. quantity

| technique | defects patched |
|---|---|
| GenProg | 49 (13.7%) |
| Par | 38 (10.6%) |
| SimFix | 68 (19.0%) |
| TRPAutoRepair | 44 (12.3%) |
| total | 106 (29.7%) |

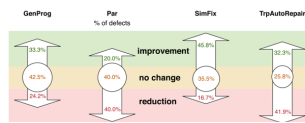**When applied to real-world Java code, APR produces patches for 10.6-19.0% of the defects**

50

### Slide 51

## Quality vs. quantity

| technique | minimum | patch quality mean | median | maximum | 100%-quality patches |
|---|---|---|---|---|---|
| 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% |

**Less than half (14-46%) of the patches are correct**

51

### Slide 52

## Does APR at least improve things a bit?

| | | | change in quality due to patch | | |
|---|---|---|---|---|---|
| technique | minimum | mean | median | maximum |
| GenProg | −30.9% | −1.7% | 0.0% | 2.6% |
| Par | −30.9% | −2.8% | 0.0% | 1.5% |
| SimFix | −24.9% | 0.2% | 0.0% | 35.0% |
| TrpAutoRepair | −30.9% | −2.1% | 0.0% | 3.8% |

52

### Slide 53

## So is there no hope?

- SearchRepair, a brand new technique, reduces overfitting to 97.2%.
- Most SearchRepair repairs pass 100% of the held-out test suite.
  (Select few poor repairs drop the overall rate.)

Read more about SearchRepair:

http://people.cs.umass.edu/~brun/pubs/pubs/Ke15ase.pdf

53

### Slide 54

Takeaway: Tests are an imperfect oracle, so APR suffers, producing low-quality patches.

Can we find a domain with better oracles?

54