# Lecture 21: Program Repair

17-355/17-655/17-819: Program Analysis

Rohan Padhye and Jonathan Aldrich
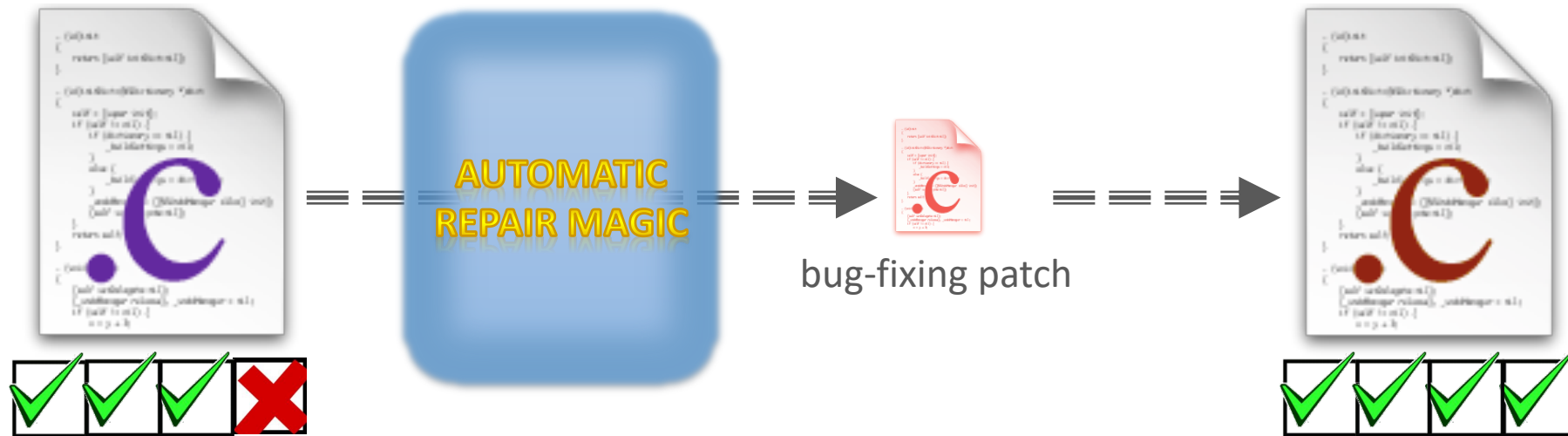
(material heavily borrowed from Claire Le Goues)

April 20, 2021

**Carnegie Mellon University**
School of Computer Science

# We've spent a lot of time on *finding* bugs

- What about *fixing* them?

- Problem: Given a program and an indication of a bug, find a patch for that program to fix that bug.
  - Both static and dynamic techniques have been used to "indicate" bugs.
  - The bulk of repair research is *dynamic*, or uses tests.
  - (We'll talk about static briefly, and again later.)

# Automatic Program Repair



bug-fixing patch

institute for SOFTWARE RESEARCH

**Carnegie Mellon University**
School of Computer Science

# Bug fixing: the 30000-foot view

1. Localize the bug.
   o And perform additional analysis
2. Create/combine fix possibilities into 1+ possible patches.
3. Validate candidate patches.

Fault localization

Tests.

isr institute for SOFTWARE RESEARCH

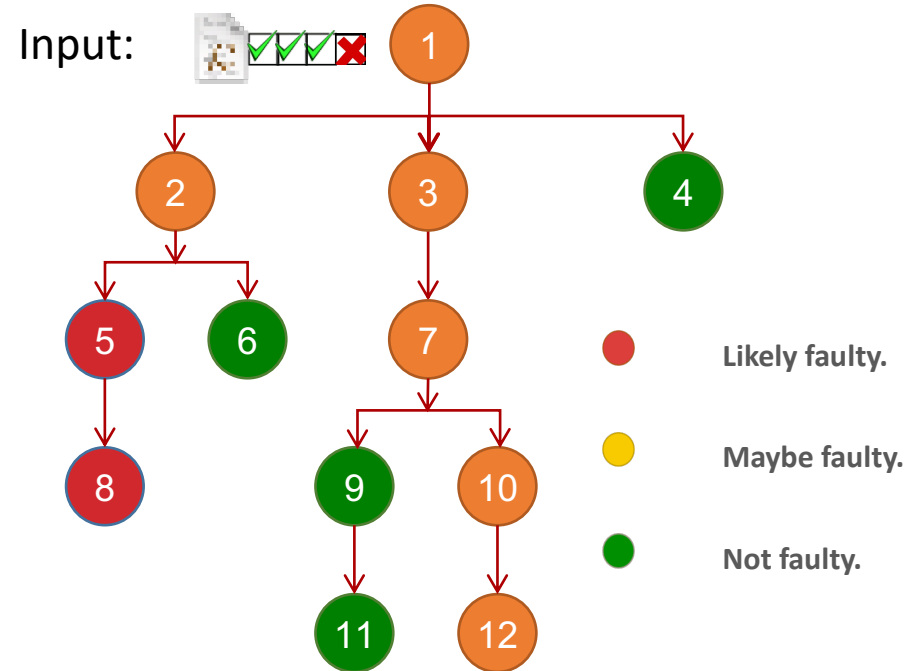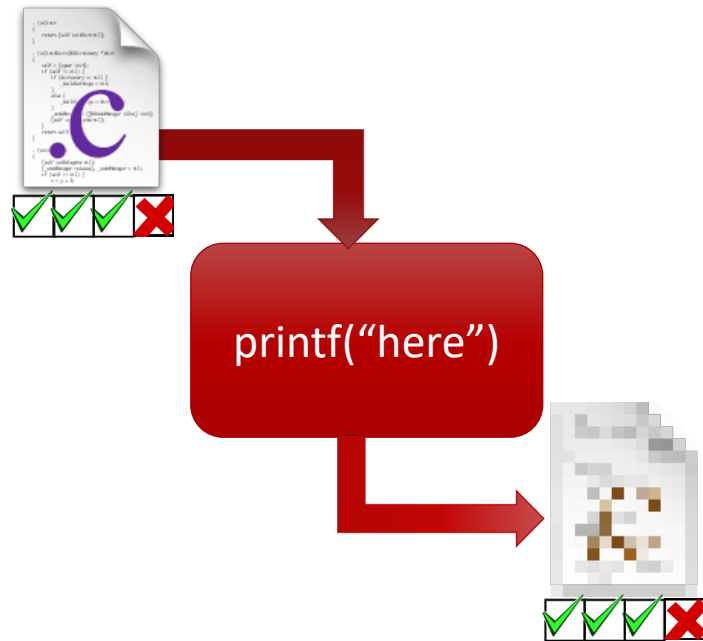**Carnegie Mellon University**
School of Computer Science

# Fault Localization

- Given: set of test cases, some of which fail

- To find: part of the code that's causing the failure
  - (which needs to be fixed)

- How is this done manually?
  - Printf("here")

# Spectrum-Based Fault Localization

Automatically ranks potentially buggy program pieces based on test case behavior.

# GenProg: Repair with Evolutionary Computation

1. Localize the bug.
   o And perform additional analysis
2. Create/combine fix possibilities into 1+ possible patches.
3. Validate candidate patches.

"GenProg: A generic method for automatic software repair" by Le Goues et al. IEEE TSE (2011)

Localize to C statements

Genetic programming

# GenProg: Repair with Evolutionary Computation

Biased, random search for AST-level edits to a program that fixes a given bug without breaking any previously-passing tests.
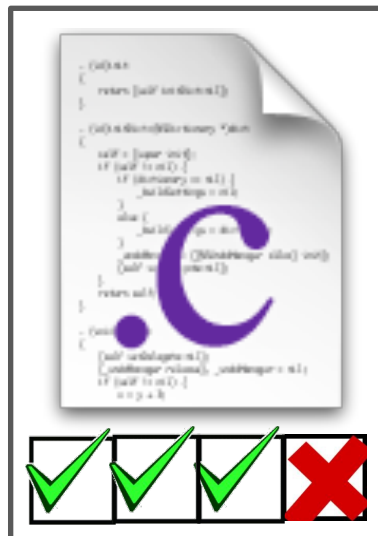
isr institute for SOFTWARE RESEARCH

**Carnegie Mellon University**
School of Computer Science

# Genetic Programming
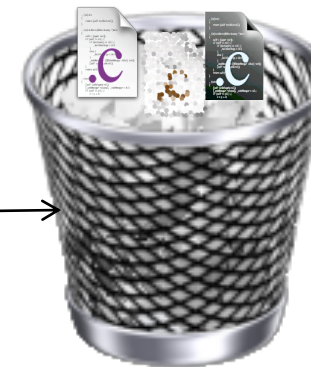
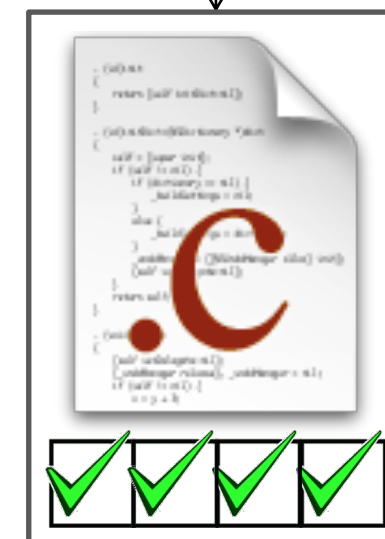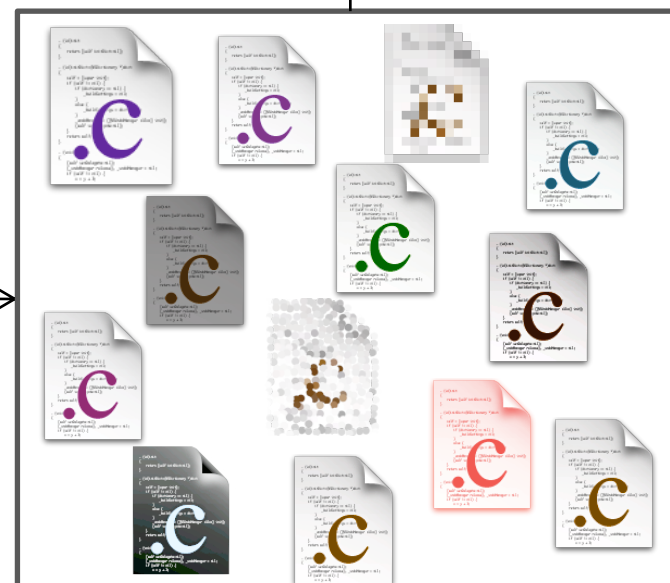**The application of evolutionary or genetic algorithms to program source code.**
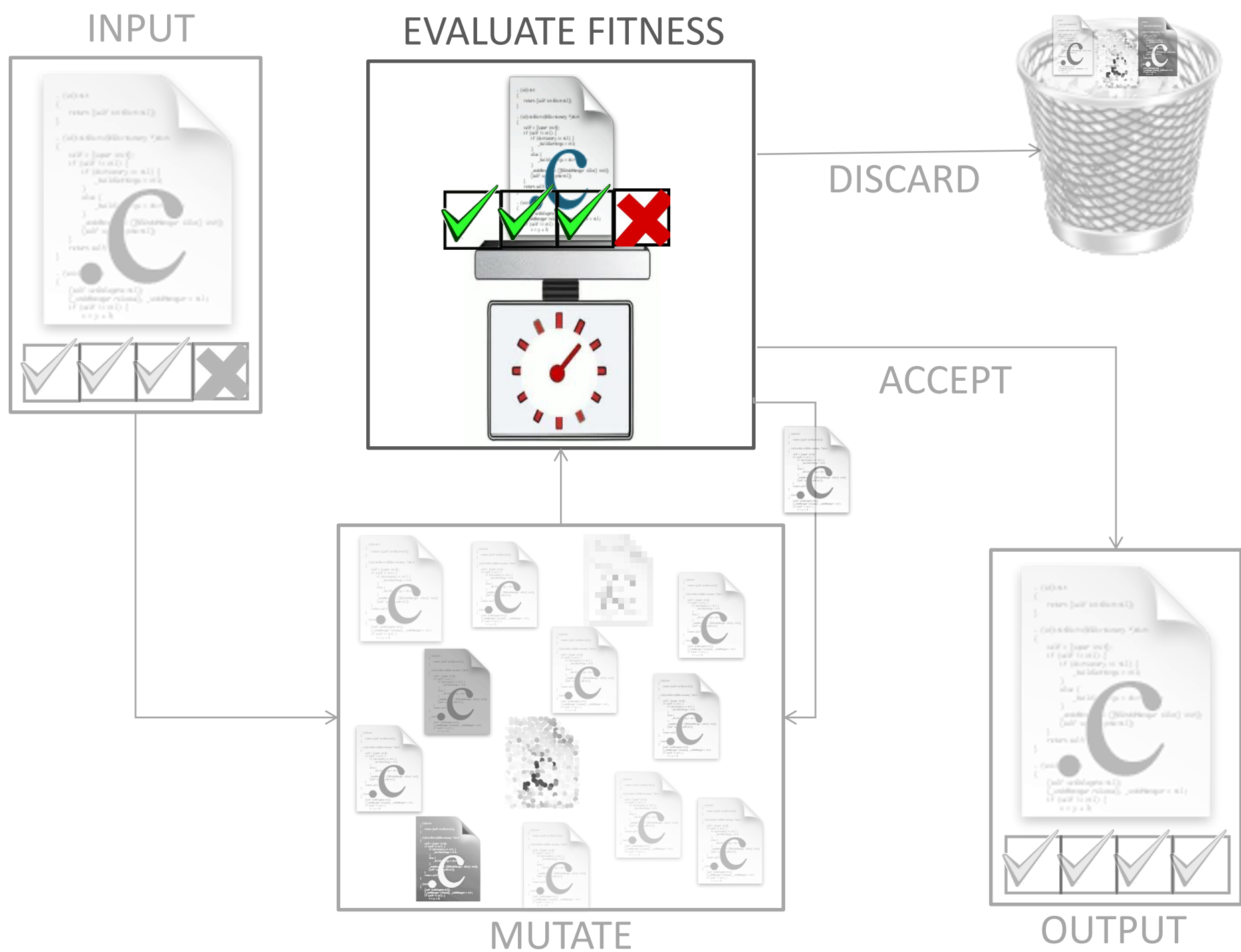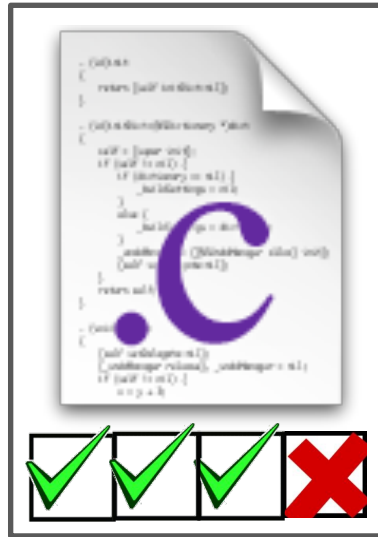
INPUT

EVALUATE FITNESS

DISCARD

ACCEPT

MUTATE

OUTPUT

INPUT

EVALUATE FITNESS

DISCARD

ACCEPT

MUTATE

OUTPUT

INPUT

EVALUATE FITNESS
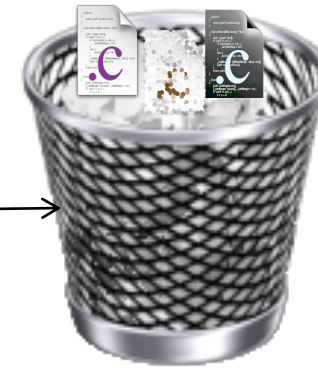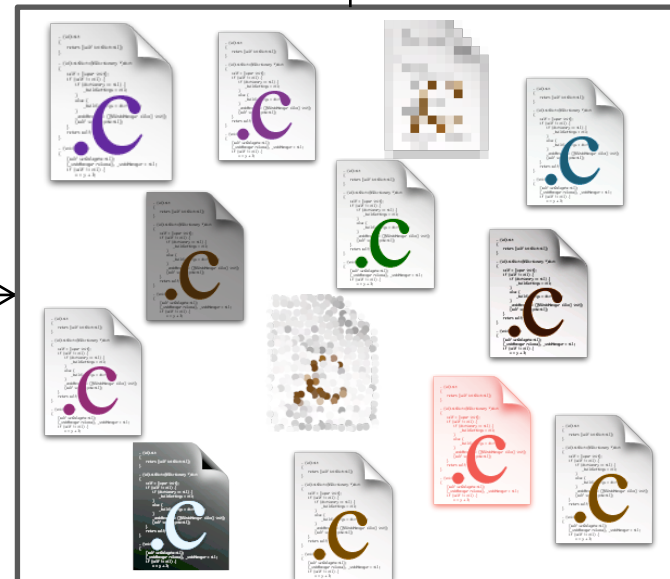
DISCARD

ACCEPT

MUTATE

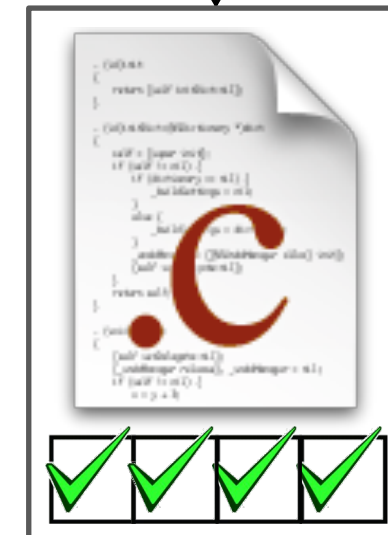OUTPUT

12

INPUT

EVALUATE FITNESS

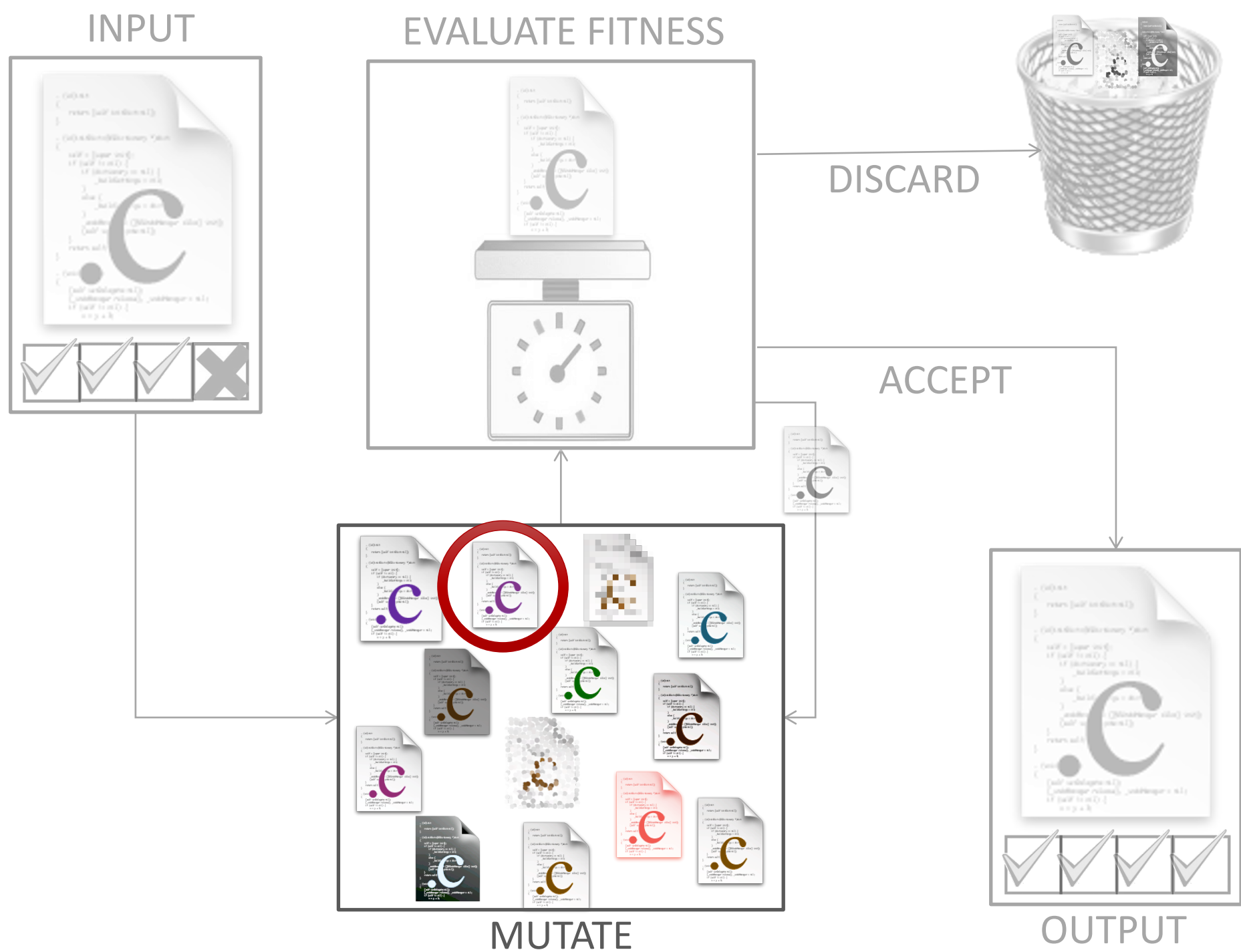DISCARD

ACCEPT

MUTATE

OUTPUT

# An individual is a candidate patch/set of changes to the input program.

- A patch is a series of *statement-level* edits:
  - delete X
  - replace X with Y
  - insert Y after X.

- Replace/insert: pick Y from somewhere else in the program.

- To mutate an individual, add new random edits to a given (possibly empty) patch.
  - (Where? Right: fault localization!)

>

```c
1  void gcd(int a, int b) {
2    if (a == 0) {
3      printf("%d", b);
4    }
5    while (b > 0) {
6      if (a > b)
7        a = a – b;
8      else
9        b = b – a;
10   }
11   printf("%d", a);
12   return;
13 }
```
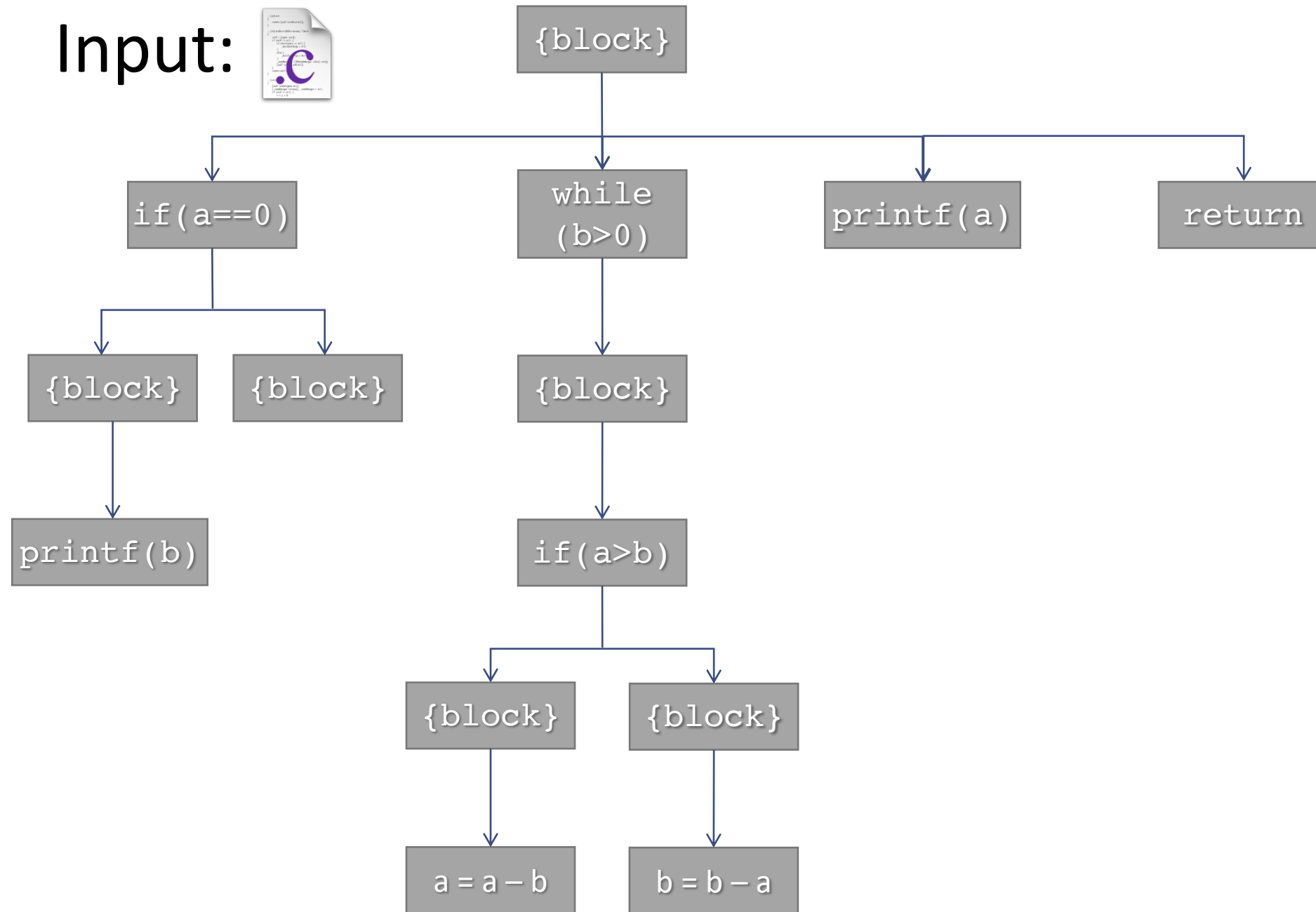
```
> gcd(4,2)
> 2
>
> gcd(1071,1029)
> 21
>
> gcd(0,55)
> 55
```

(looping forever)

```c
1  void gcd(int a, int b) {
2      if (a == 0) {
3          printf("%d", b);
4      }
5      while (b > 0) {
6          if (a > b)
7              a = a – b;
8          else
9              b = b – a;
10     }
11     printf("%d", a);
12     return;
13 }
```

!

Input: 

```
{block}
```

```
if(a==0)        while        printf(a)        return
                (b>0)
```

```
{block}    {block}        {block}
```

```
printf(b)                if(a>b)
```

```
{block}    {block}
```

$$a = a - b$$            $$b = b - a$$

Input: ✅✅✅❌



An **edit** is:

- Insert statement X after statement Y
- Replace statement X with statement Y
- Delete statement X

Input: ✓ ✓ ✓ ✗

{block}

if(a==0)

while
(b>0)

printf(a)

return

{block}

{block}

{block}

printf(b)

if(a>b)

{block}

{block}

a = a − b

b = b − a

An **edit** is:

- **Insert statement X after statement Y**
- Replace statement X with statement Y
- Delete statement X

Input: ✓✓✓✗

{block}

if(a==0)  while(b>0)  printf(a)  return

{block}  {block}  {block}

printf(b)  if(a>b)

**An edit is:**

- **Insert statement X after statement Y**
- Replace statement X with statement Y
- Delete statement X

{block}  {block}

a = a − b  b = b − a

21

Input: ☑ ☑ ☑ ☑

{block}

if(a==0)  while (b>0)  printf(a)  return

{block}  {block}  {block}

printf(b)  if(a>b)

return

**An edit is:**

- **Insert statement X after statement Y**
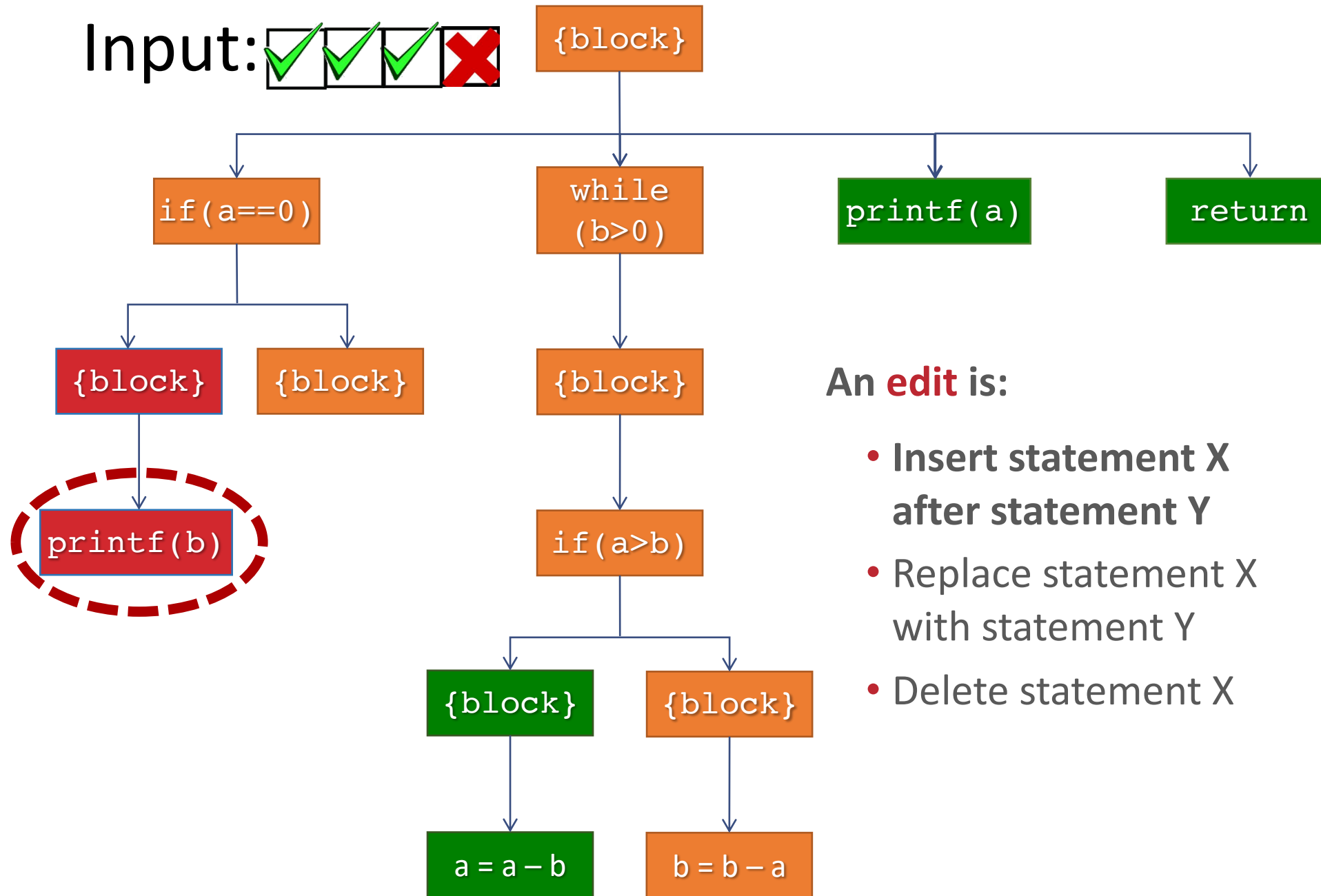- Replace statement X with statement Y
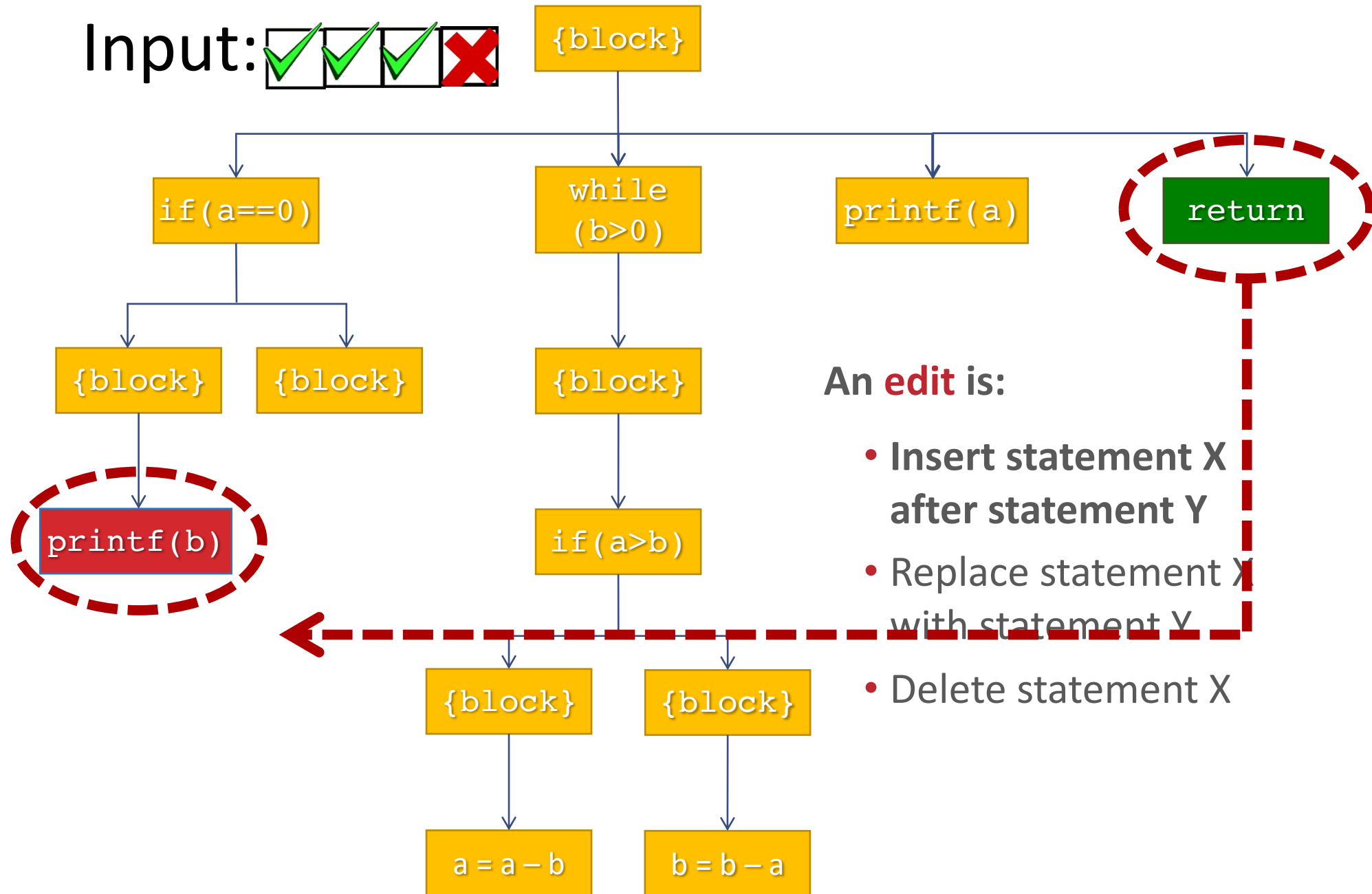- Delete statement X

{block}  {block}

a = a − b  b = b − a

22
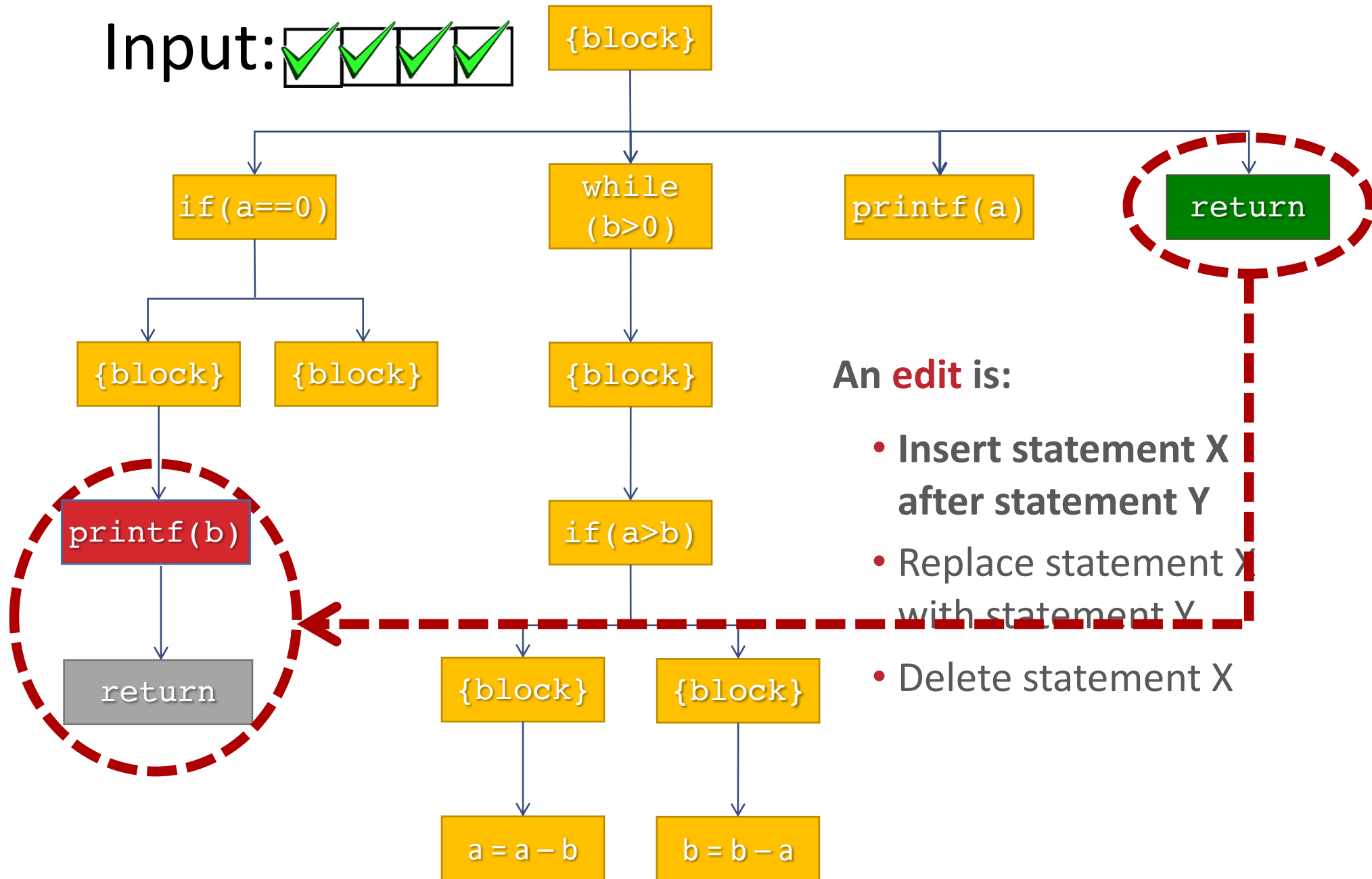
# Wait.. Isn't this expensive?

## A Systematic Study of Automated Program Repair:
## Fixing 55 out of 105 Bugs for $8 Each

| Claire Le Goues | Michael Dewey-Vogt | Stephanie Forrest | Westley Weimer |
|---|---|---|---|
| Computer Science Department | | Computer Science Department | Computer Science Department |
| University of Virginia | | University of New Mexico | University of Virginia |
| Charlottesville, VA | | Albuquerque, NM | Charlottesville, VA |
| {legoues,mkd5m}@cs.virginia.edu | | forrest@cs.unm.edu | weimer@cs.virginia.edu |

*Abstract*—There are more bugs in real-world programs than human programmers can realistically address. This paper evaluates two research questions: "What fraction of bugs can be repaired automatically?" and "How much does it cost to repair a bug automatically?" In previous work, we presented *GenProg*, which uses genetic programming to repair defects in off-the-shelf C programs. To answer these questions, we: (1) propose novel algorithmic improvements to GenProg that allow it to scale to large programs and find repairs 68% more often, (2) exploit GenProg's inherent parallelism using cloud computing resources to provide grounded, human-competitive cost measurements, and (3) generate a large, indicative benchmark set to use for systematic evaluations. We evaluate GenProg on 105 defects from 8 open-source programs totaling 5.1 million lines of code and involving 10,193 test cases. GenProg automatically repairs 55 of those 105 defects. To our knowledge, this evaluation is the largest available of its kind, and is often two orders of magnitude larger than previous work in terms of code or test suite size or defect count. Public cloud computing prices allow our 105 runs to be reproduced for $403; a successful repair completes in 96 minutes and costs $7.32, on average.

*Keywords*-genetic programming; automated program repair; cloud computing

patch overflow and illegal control-flow transfer vulnerabilities; AutoFix-E [9], which can repair programs annotated with design-by-contract pre- and post-conditions; and AFix [10], which can repair single-variable atomicity violations. In previous work, we introduced *GenProg* [11], [12], [13], [14], a general method that uses genetic programming (GP) to repair a wide range of defect types in legacy software (e.g., infinite loops, buffer overruns, segfaults, integer overflows, incorrect output, format string attacks) without requiring *a priori* knowledge, specialization, or specifications. GenProg searches for a repair that retains required functionality by constructing variant programs through computational analogs of biological processes.

The goal of this paper is to evaluate dual research questions: "What fraction of bugs can GenProg repair?" and "How much does it cost to repair a bug with GenProg?" We combine three important insights to answer these questions. Our key algorithmic insight is to represent candidate repairs as patches [15], rather than as abstract syntax trees. These changes were critical to GenProg's scalability to millions of lines of code, an essential component of our evaluation.

23

# "Real world" applications

**"Facebook, Inc"**

**"one widely-studied [repair] approach uses software testing to guide the repair process, as typified by GenProg."**

**"Results from repair applied to 6 multi-million line systems."**

# Can GenProg fix this?

- The checksum program should:
  - Take a single-line string as input.
  - Sum the integer codes of the characters, excluding the newline, modulo 64, plus the code for the space character.

- Buggy student assignment →

```
1. // …
2. while (next != '\n')
3. {
4.     scanf("%c", &next);
5.     sum += next;
6. }
7. sum = sum % 64 + 22;
       um;
```

**Incorrectly includes the newline in the sum.**

**Wrong value: the ASCII value of space is 32, not 22.**

# Voila!

- The checksum program should:
  o Take a single-line string as input.
  o Sum the integer codes of the characters in the string, modulo 64, plus the code for the space character.

- GenProg fix with new representation →

```
1. // …
2. while (next != '\n')
3. {
+    FIXME scanf("%c", &next);
4.    sum += next;
+ if (next == '\n')
+      break;
4. }
5. sum = sum % 64 + 22;
+    sum += next;
8. return sum;
```

# Semantics-based repair

1. Localize the bug.
   o And perform additional analysis
2. Create/combine fix possibilities into 1+ possible patches.
3. Validate candidate patch.

Same idea, but localizing to *expressions*.

RHS of assignments, conditionals.

"SemFix: Program Repair via Semantic Analysis" by Nguyen et al. ICSE 2013

"Angelix: Scalable Multiline Program Patch Synthesis via Symbolic Analysis" by Mechtaev et al. ICSE 2016

```
1   int is_upward( int inhibit, int up_sep, int down_sep){
2           int bias;
3           if (inhibit)
4               bias = down_sep;
5           else  bias = up_sep ;
6           if (bias > down_sep)
7                   return 1;
8           else return 0;
9   }
```

(Slides by Abhik Roychoudhury)

```
1  int is_upward( int inhibit, int up_sep, int down_sep){
2       int bias;
3       if (inhibit)
4            bias = down_sep;
5       else  bias = up_sep ;
6       if (bias > down_sep)
7            return 1;
8       else return 0;
9  }
```

| inhibit | up_sep | down_sep | Observed output | Expected Output | Result |
|---------|--------|----------|-----------------|-----------------|--------|
| 1 | 0 | 100 | 0 | 0 | pass |
| 1 | 11 | 110 | 0 | 1 | fail |
| 0 | 100 | 50 | 1 | 1 | pass |
| 1 | -20 | 60 | 0 | 1 | fail |
| 0 | 0 | 10 | 0 | 0 | pass |

```
1  int is_upward( int inhibit, int up_sep, int down_sep){
2        int bias;
3        if (inhibit)
4             bias = down_sep;
5        else  bias = up_sep ;
6        if (bias > down_sep)
7             return 1;
8        else return 0;
9  }
```

| inhibit | up_sep | down_sep | Observed output | Expected Output | Result |
|---|---|---|---|---|---|
| 1 | 0 | 100 | 0 | 0 | pass |
| 1 | 11 | 110 | 0 | 1 | fail |
| 0 | 100 | 50 | 1 | 1 | pass |
| 1 | -20 | 60 | 0 | 1 | fail |
| 0 | 0 | 10 | 0 | 0 | pass |

```
1   int is_upward( int inhibit, int up_sep, int down_sep){
2         int bias;
3         if (inhibit)
4              bias = down_sep;   // bias= up_sep + 100
5         else  bias = up_sep ;
6         if (bias > down_sep)
7              return 1;
8         else return 0;
9   }
```

| inhibit | up_sep | down_sep | Observed output | Expected Output | Result |
|---------|--------|----------|-----------------|-----------------|--------|
| 1 | 0 | 100 | 0 | 0 | pass |
| 1 | 11 | 110 | 0 | 1 | fail |
| 0 | 100 | 50 | 1 | 1 | pass |
| 1 | -20 | 60 | 0 | 1 | fail |
| 0 | 0 | 10 | 0 | 0 | pass |

# Angelix

1. Localize the bug.
   o And perform additional...
2. Create/combine fix possibilities
   into 1+ possible patches.
3. Validate candidate patch...

*Concolic execution* to find expression values that would make the test pass.

*Program synthesis* to construct replacement code that produces those values.

# An expression's *angelic value* is the value that would make a given test case pass.

- This value is set "arbitrarily", by which we mean symbolically.

- You can *solve* for this value if you have:
  o the test case's expected input/output.
  o the path condition controlling its execution.

- Concolic execution (remember me?):
  o Start executing the test concretely, and then switch to symbolic execution when the angelic value starts to matter.

```
1  int is_upward( int inhibit, int up_sep, int down_sep){
2       int bias;
3       if (inhibit)
4            bias = down_sep;
5       else  bias = up_sep ;
6       if (bias > down_sep)
7            return 1;
8       else return 0;
9  }
```

| inhibit | up_sep | down_sep | Observed output | Expected Output | Result |
|---------|--------|----------|-----------------|-----------------|--------|
| 1 | 0 | 100 | 0 | 0 | pass |
| 1 | 11 | 110 | 0 | 1 | fail |
| 0 | 100 | 50 | 1 | 1 | pass |
| 1 | -20 | 60 | 0 | 1 | fail |
| 0 | 0 | 10 | 0 | 0 | pass |

```
1   int is_upward( int inhibit, int up_sep, int down_sep){
2        int bias;
3        if (inhibit)
4            bias = α;
5        else  bias = up_sep ;
6        if (bias > down_sep)
7            return 1;
8        else return 0;
9   }
```

| inhibit | up_sep | down_sep | Observed output | Expected Output | Result |
|---------|--------|----------|-----------------|-----------------|--------|
| 1 | 11 | 110 | 0 | 1 | fail |

inhibit = 1, up_sep = 11, down_sep = 110
bias = α,  PC = true

Line 4

Line 7

inhibit = 1, up_sep = 11, down_sep = 110
bias = α, PC= α > 110

Line 8

inhibit = 1, up_sep = 11, down_sep = 110
bias = α, PC= α > 110

35

```
1   int is_upward( int inhibit, int up_sep, int down_sep){
2         int bias;
3         if (inhibit)
4             bias = α;
5         else  bias = up_sep ;
6         if (bias > down_sep)
7               return 1;
8         else return 0;
9   }
```

**Exercise**: Generate constraints for all other test cases

| inhibit | up_sep | down_sep | Observed output | Expected Output | Result | Constraint |
|---------|--------|----------|-----------------|-----------------|--------|------------|
| 1 | 0 | 100 | 0 | 0 | pass | |
| 1 | 11 | 110 | 0 | 1 | fail | f(1,11,110) > 110 |
| 0 | 100 | 50 | 1 | 1 | pass | |
| 1 | -20 | 60 | 0 | 1 | fail | |
| 0 | 0 | 10 | 0 | 0 | pass | |

# Collect all of the constraints!

- Accumulated constraints over all test cases:

$$f(1,11,110) > 110 \land f(1,0,100) \leq 100$$
$$\land f(1,-20,60) > 60$$

- Use **oracle guided component-based program synthesis** to construct satisfying *f*:
  - *How does this work again?*

- Generated fix
  - `f(inhibit,up_sep,down_sep) = up_sep + 100`

# Heartbleed patch

```
if (hbtype == TLS1_HB_REQUEST
    && (payload + 18) < s->s3->rrec.length) {

   ...
} else if (hbtype == TLS1_HB_RESPONSE) {
   ...
}
return 0;



if (1 + 2 + payload + 16 > s->s3->rrec.length)
   return 0;
   ...
if (hbtype == TLS1_HB_REQUEST) {
   ...
} else if (hbtype == TLS1_HB_RESPONSE) {
   ...
}
return 0;
```

Generated patch

Developer patch

# Challenges & Trade-Offs

1. Scalability

2. Expressibility of Repair

3. Patch Quality

# Open problem: What is a high quality patch?

- Understandable?

- Doesn't delete stuff?
  o Think: `assert(p) ---> ~~assert(p)~~`

- Addresses the cause, not the symptom...
  o Think: `++a[i] ---> try{++a[i]}catch(Exception e){}`

- Does the same thing the human did/would do?
  o But humans are often wrong!  And how close does it have to be?