

Automatic Patch Generation

Claire {Le~Goues}

PWLConf; September 15, 2016

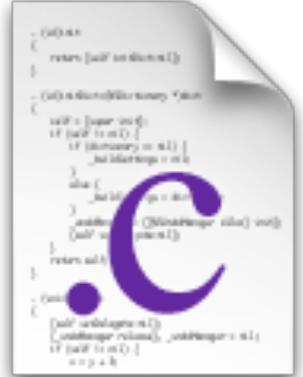
Carnegie Mellon University

isr institute for
SOFTWARE
RESEARCH

**ONCE UPON A
TIME...**



Young Claire
(developer)



Bug report from
a customer...

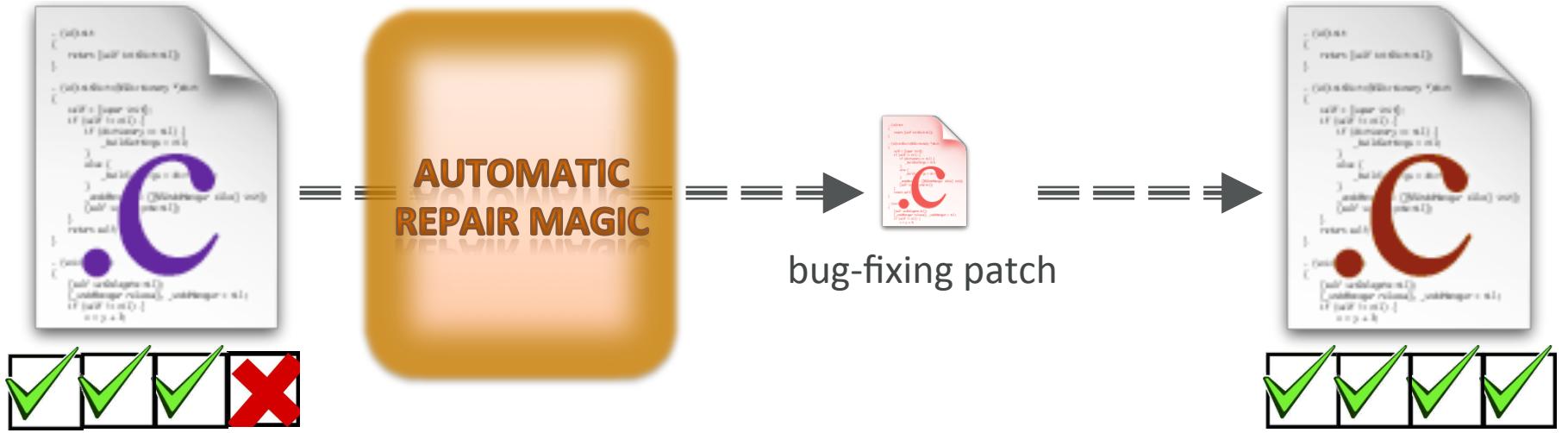
А а Аа Б б ЂЂ В в Вв Г г Ѓѓ
Д д Ђđ Е е Ѓе Ѓё Ѓё Ж ж Ђж
З з Зз И и Ии Й й й К к Ђк
Л л Ђл М м Ђми Н н Ђн О о Оо
П п Ђп Р р Ђр С с Сс Т т Ђт
У у Ђу Ф ф Ђф Х х Хх Ц ц Ђц
Ч ч Ђч Ш ш Ђши Щ щ Ђши ъ ъ
Ы ы ѿ ъ ъ Э э Ээ Ю ю Юю Я я Ђя



??!

32-64 bit transition
for Unicode encoding
of Ukrainian.

Problem: source-level defect repair



Bug fixing: the 30000-foot view

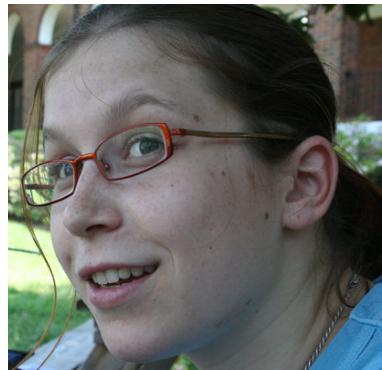
1. Localize the bug.
 - And possibly analyze it a little bit...
2. Create/combine fix possibilities into 1+ possible patches.
3. Validate candidate patches.



Bug fixing: the 30000-foot view

1. Localize the bug
 - And possibly analyze it a little bit...
2. Create/combine fix possibilities into 1+ possible patches.
3. Validate candidate patches.

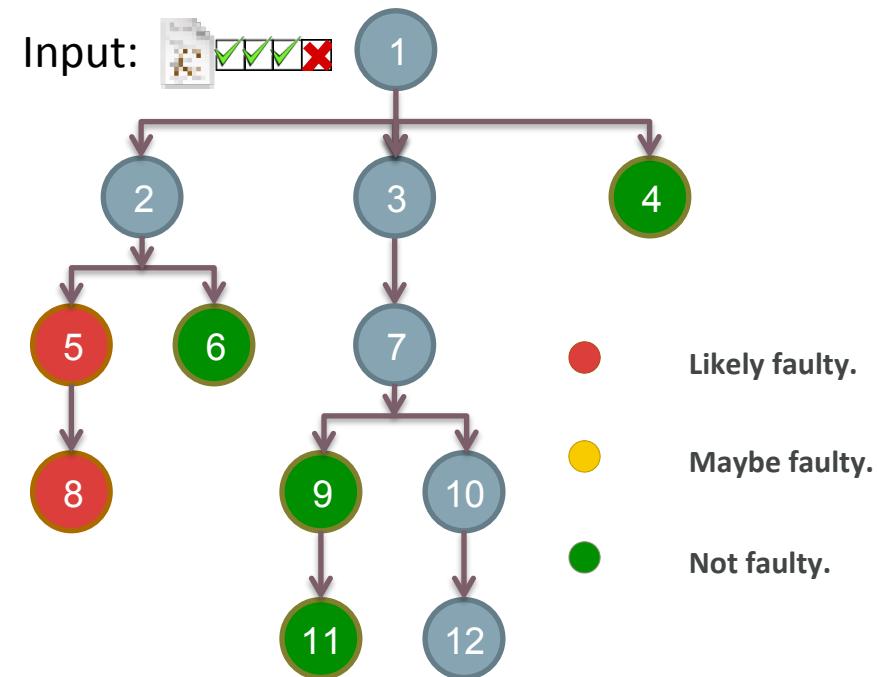
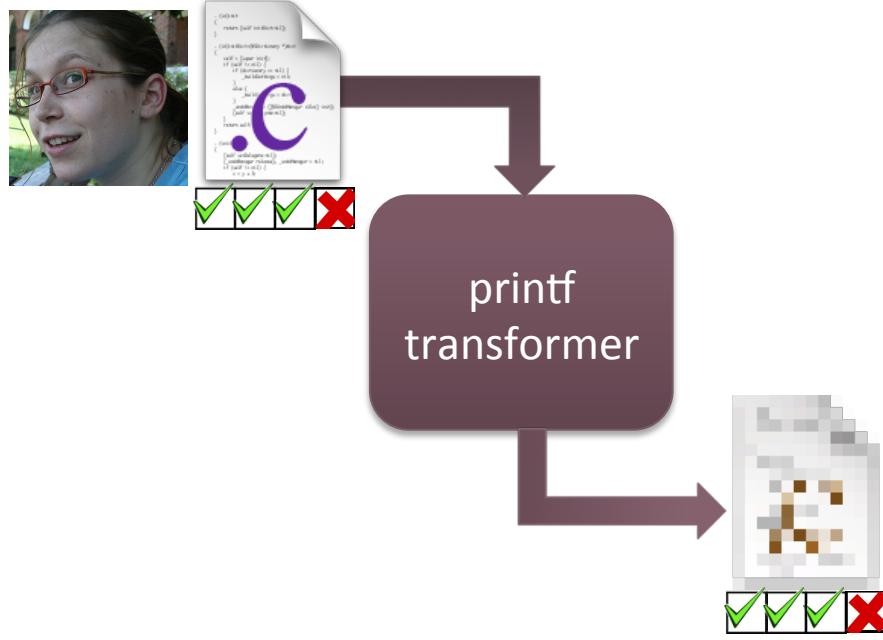




printf
transformer

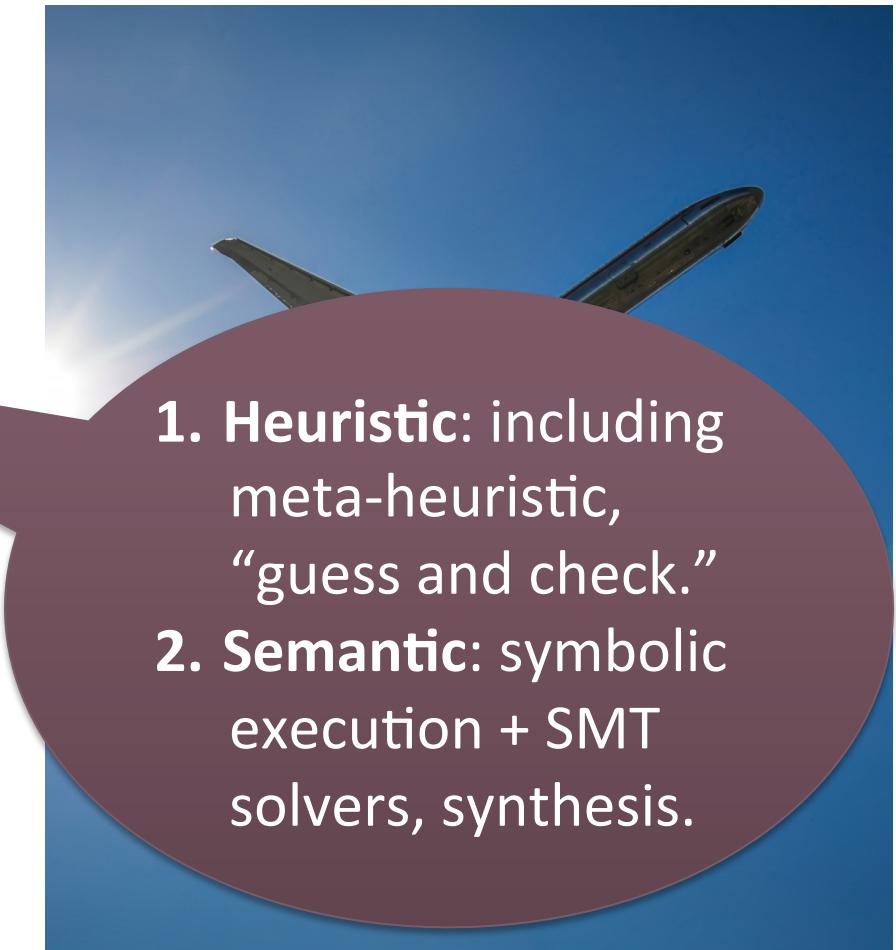


Spectrum-based fault localization automatically ranks potentially buggy program pieces based on test case behavior.



Bug fixing: the 30000-foot view

1. Localize the bug.
 - And possibly analyze it a little bit...
2. Create/combine fix possibilities into 1 or more patches.
3. Validate candidate patch.



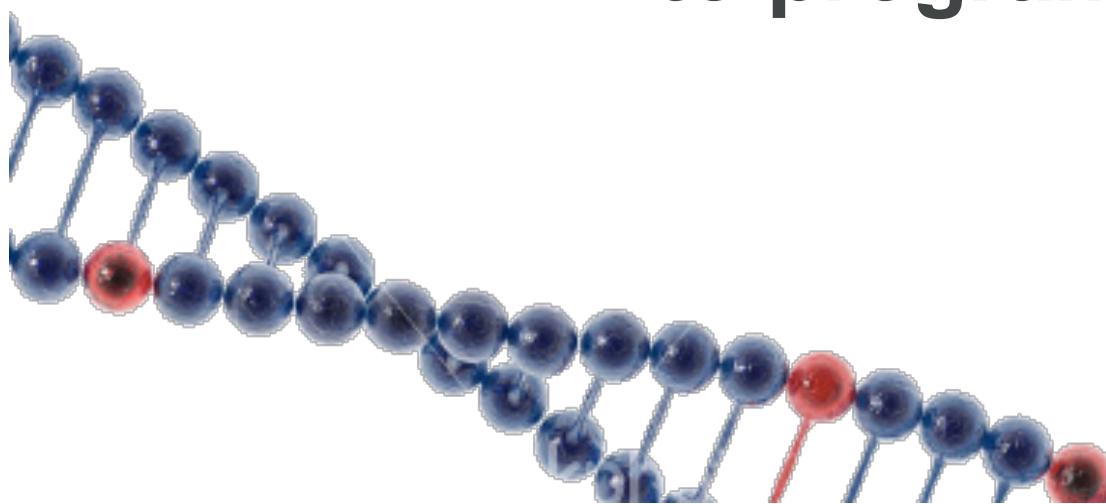
1. **Heuristic:** including meta-heuristic, “guess and check.”
2. **Semantic:** symbolic execution + SMT solvers, synthesis.

GenProg: automatic program repair using genetic programming.

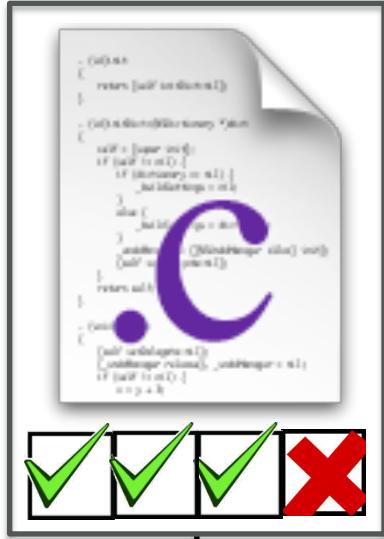
Biased, random search for a AST-level edits to a program that fixes a given bug without breaking any previously-passing tests.



Genetic programming: the application
of evolutionary or genetic algorithms
to program source code.



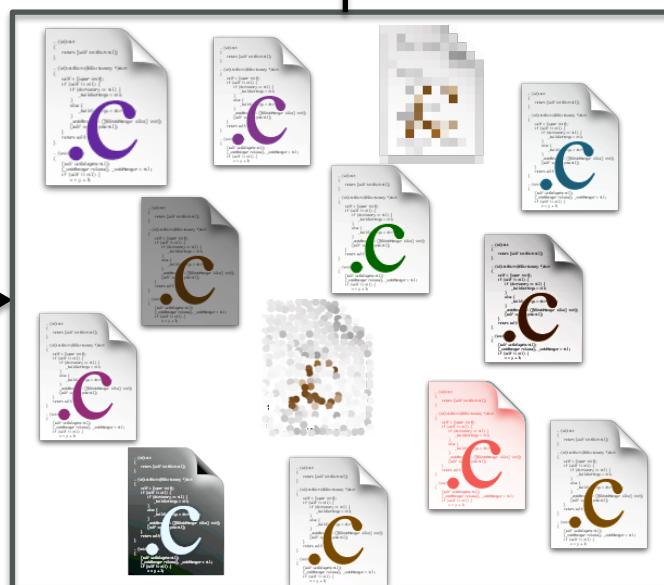
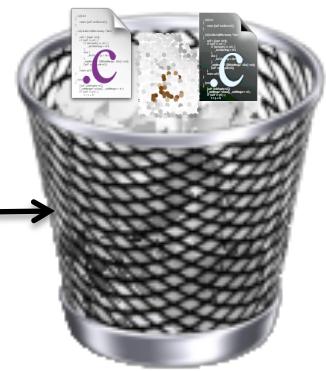
INPUT



EVALUATE FITNESS

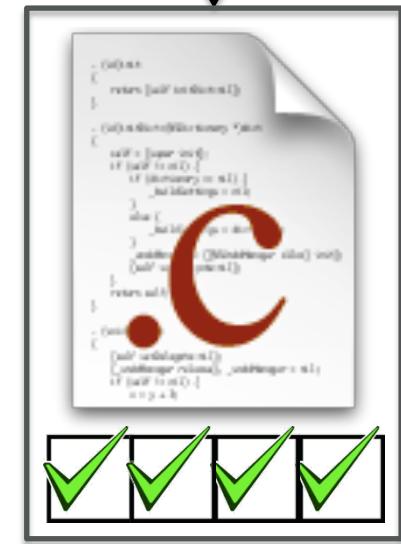


DISCARD



MUTATE

ACCEPT



OUTPUT¹⁵

GenProg: meta-heuristic approach

1. Localize the bug.
 - And possibly analyze it a little bit...
2. Create/combine fix possibilities for all possible patches.
3. Validate candidate patch.

Localize to C statements.

Use genetic programming to search for statement-level patches, reusing code from existing program.

>

```
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)

```
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 }
```

!

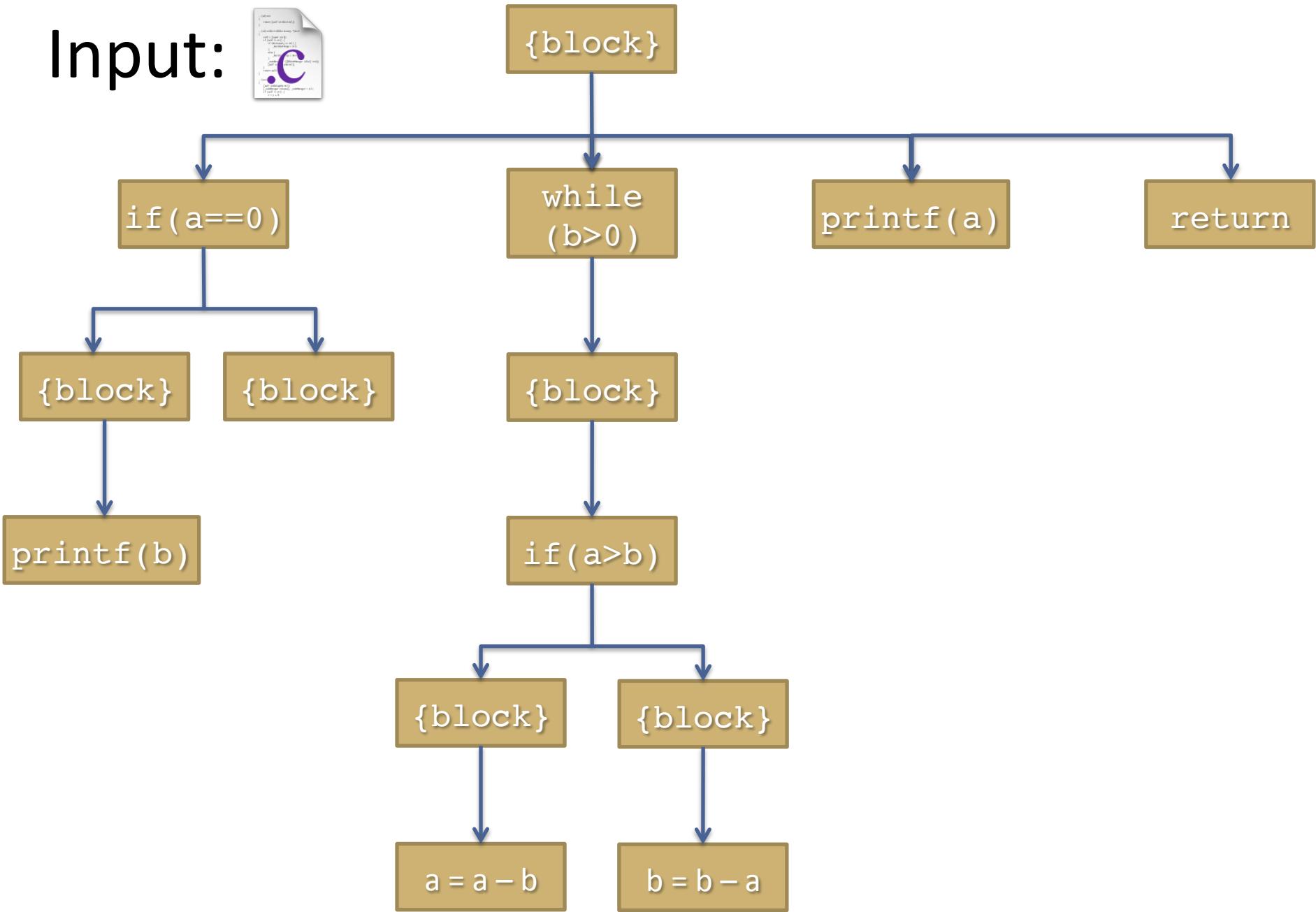
GenProg: meta-heuristic search

1. Localize the bug.
 - And possibly analyze it a little bit...
2. Create/combine fix possibilities for possible patches.
3. Validate candidate patch.

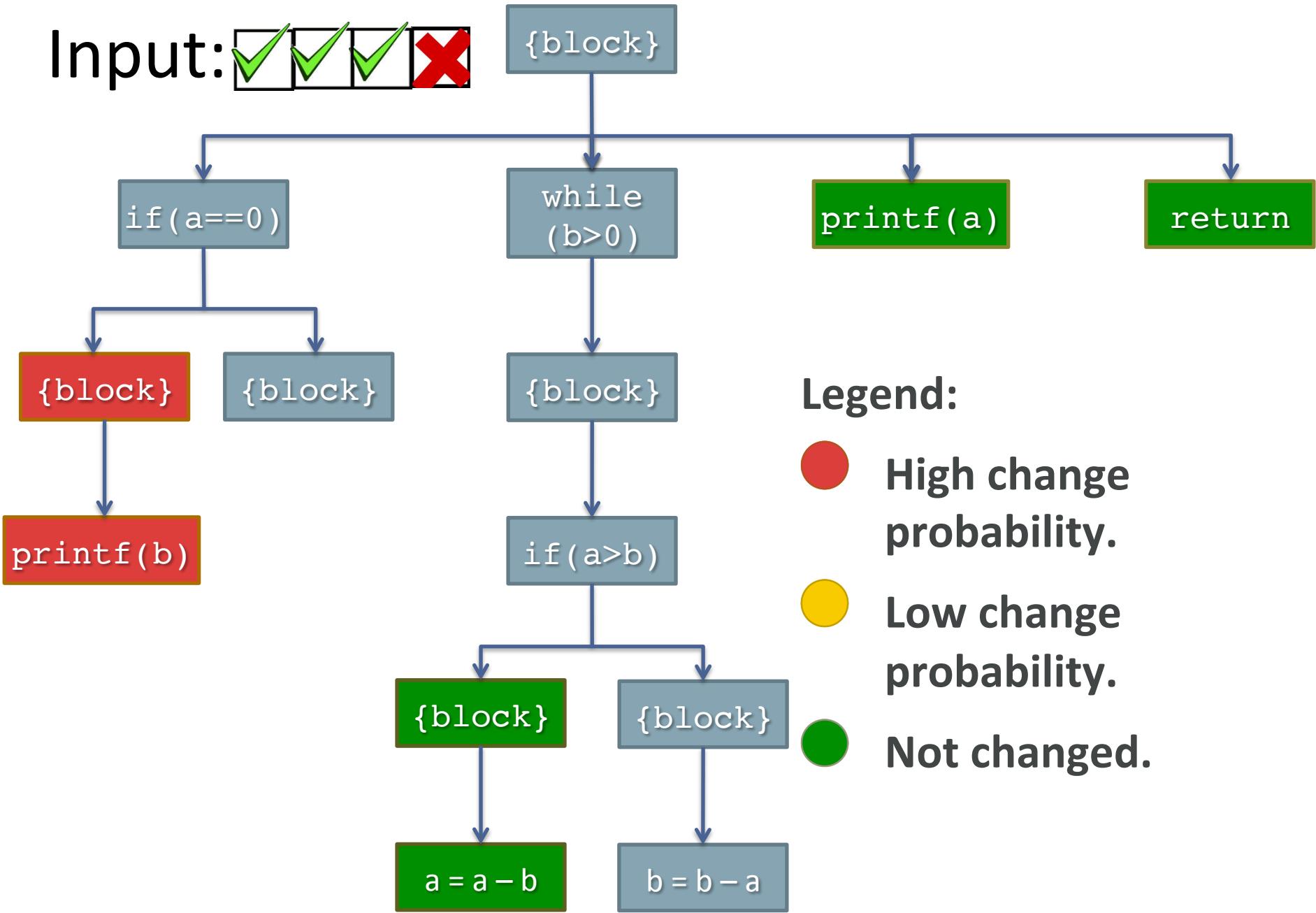
Localize to C statements.

Use genetic programming to search for statement-level patches, reusing code from existing program.

Input:



Input: 



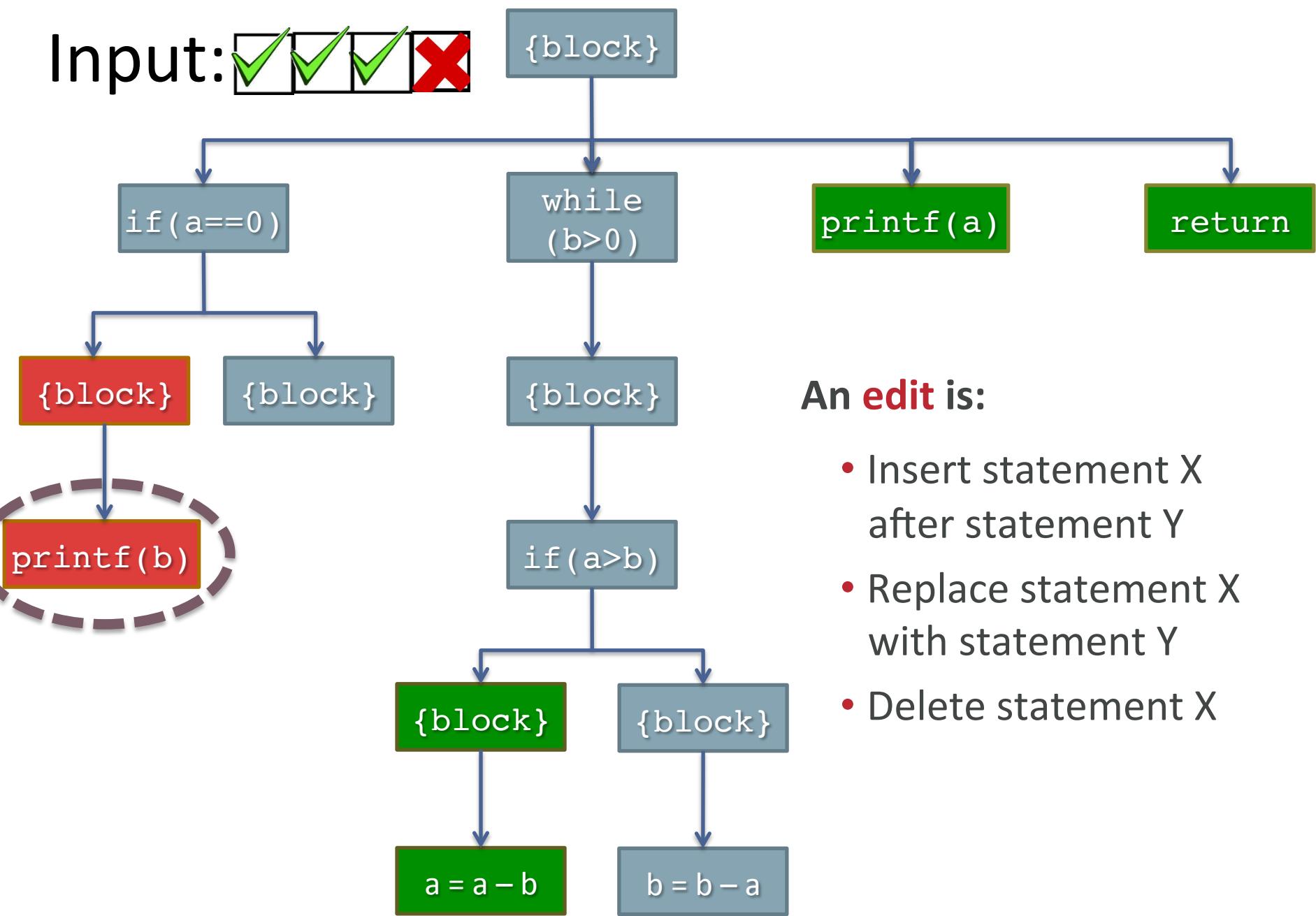
Legend:

- High change probability.
- Low change probability.
- Not changed.

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!)

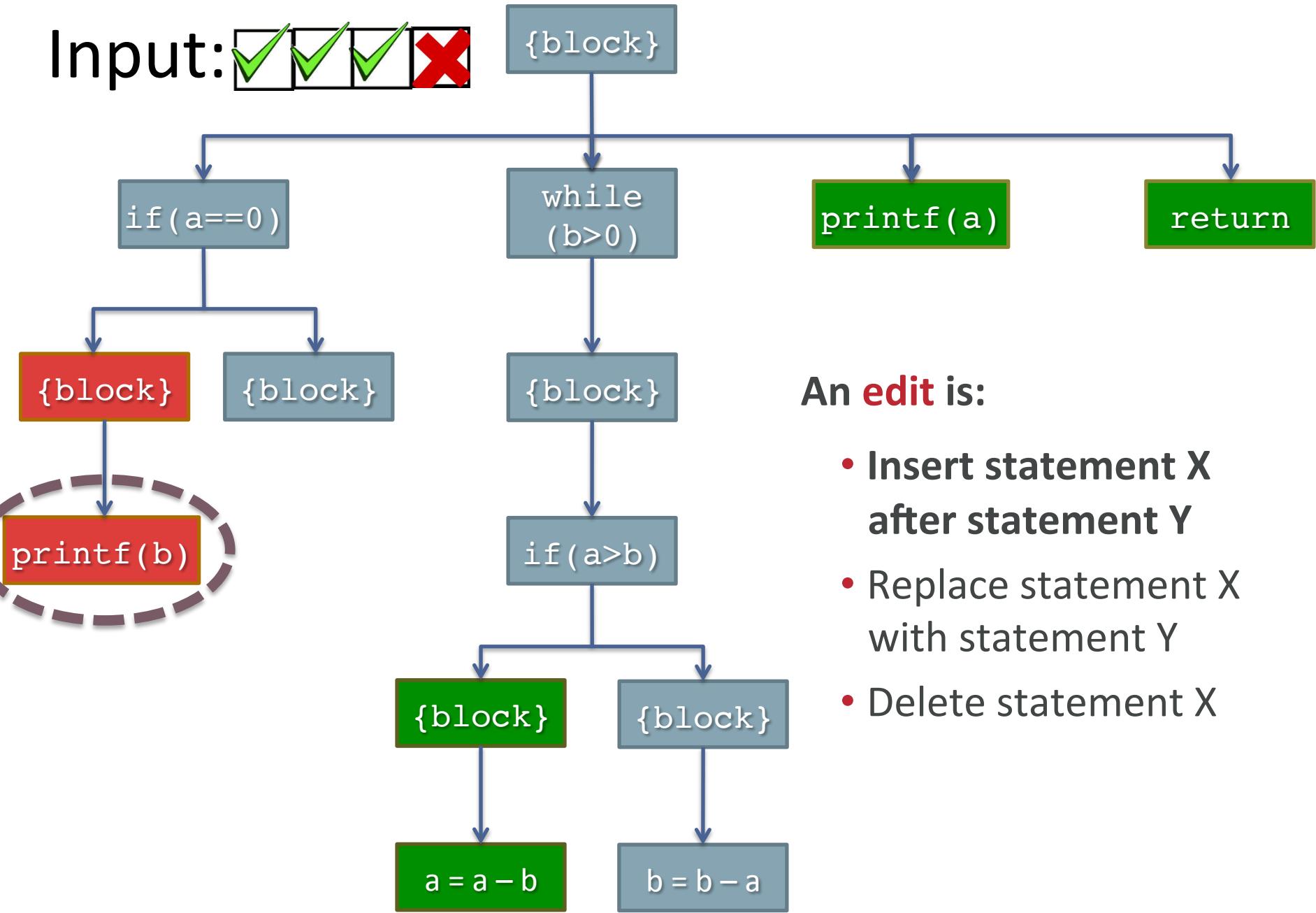
Input: 



An **edit** is:

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

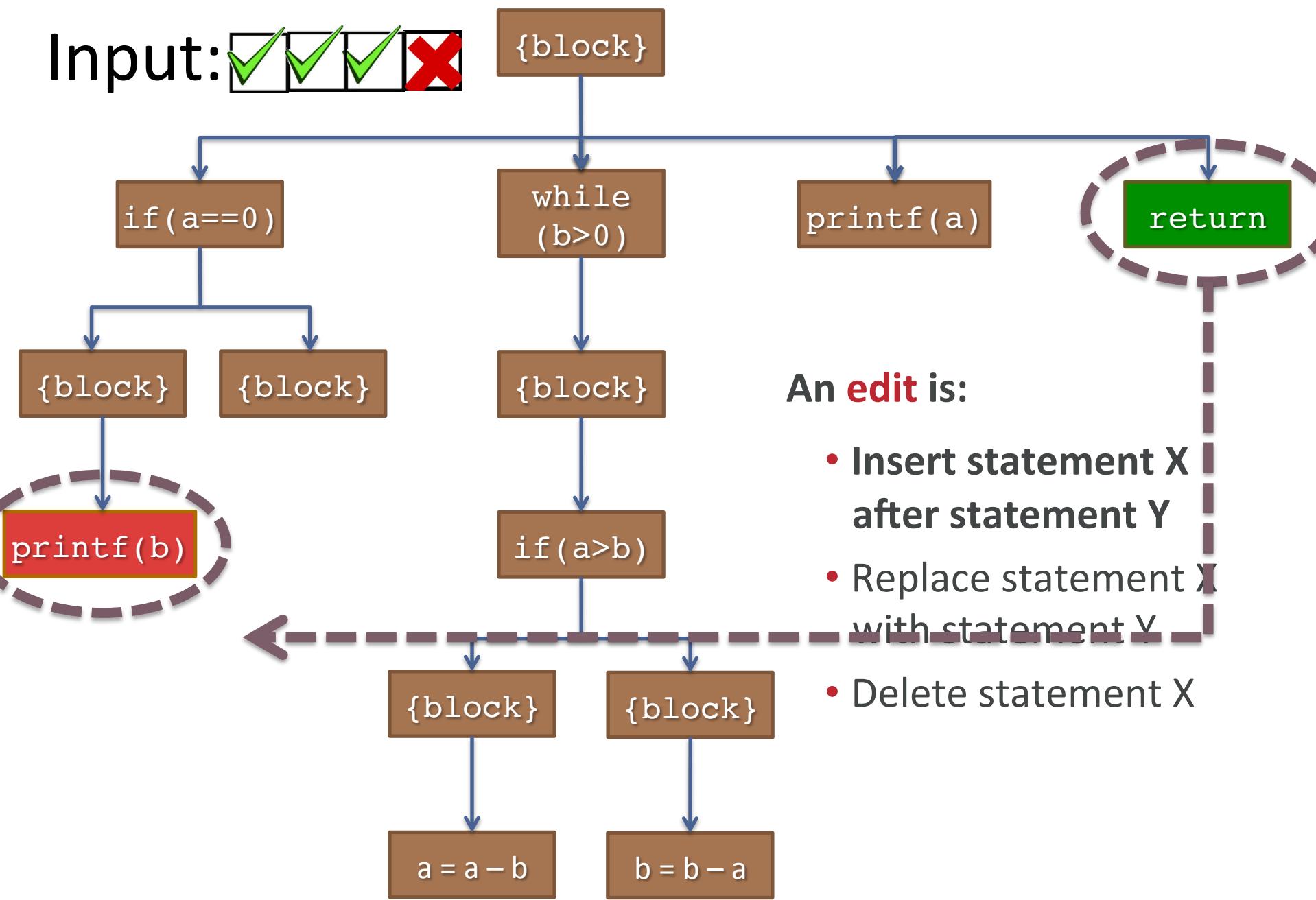
Input: 



An **edit** is:

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

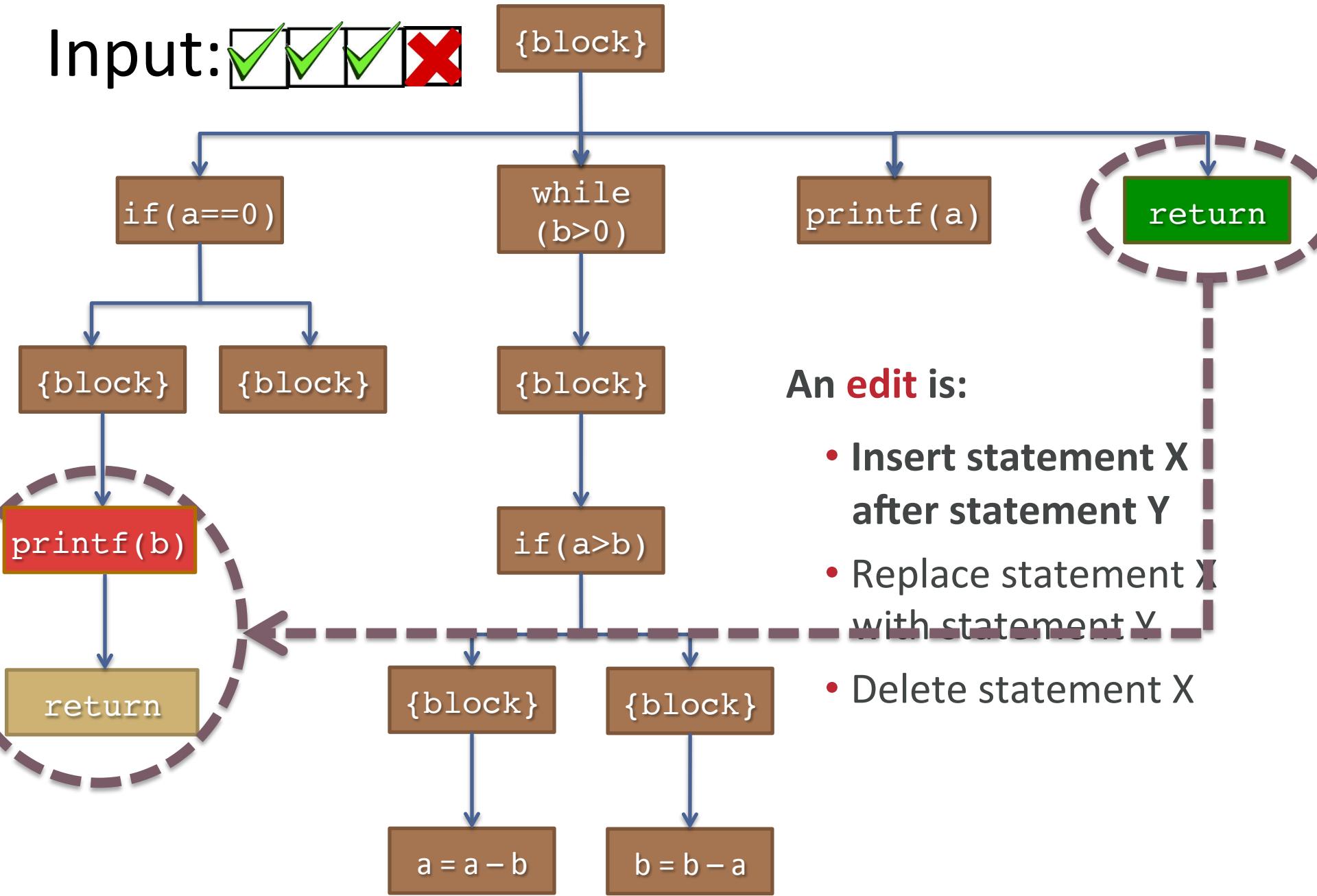
Input: 



An **edit** is:

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

Input: 



An **edit** is:

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

What about Angelic?

1. Localize the bug.
 - And possibly analyze it a little bit...
2. Create/combine fix possibilities into 1+ possible patches.
3. Validate candidate patch.

Same idea, but localizing to *expressions*.

RHS of assignments, conditionals.

```
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 }
```

Tremendous gratitude to Abhik Roychoudhury for sharing slides with me as starting material for this talk.

```

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

What about Angelix?

1. Localize the bug.
 - And possibly, a little bit...
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:
 - the test case’s expected input/output.
 - the path condition controlling its execution.
- Path condition: the set of conditions that controlled a particular execution.
 - 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; // 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

```

1 int is_upward( int inhibit, int up_sep, int down_sep) {
2     int bias;
3     if (inhibit)
4         bias =  $\alpha$ ; // 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	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= $\alpha > 110$

Line 8

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

What should it have been?

inhibit == 1

up_sep == 11

down_sep == 110

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

Symbolic Execution

$f(1,11,110) > 110$



Collect all of the constraints!

- Accumulated constraints over all test cases:

$$\begin{aligned} f(1,11,110) &> 110 \wedge f(1,0,100) \leq 100 \\ \wedge f(1,-20,60) &> 60 \end{aligned}$$

- Use oracle guided component-based program synthesis to construct satisfying f :
 - Fix a set of operators (component-based).
 - Synthesize code that only uses those operators and satisfies the constraints (oracle guided).
- Generated fix
 - $f(\text{inhibit}, \text{up_sep}, \text{down_sep}) = \text{up_sep} + 100$



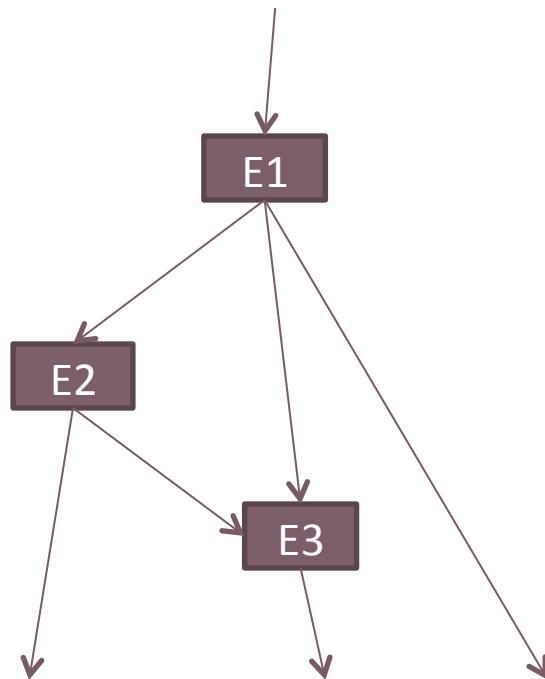
**(Legitimately interesting encoding
of synthesis problem elided for
(dubious) brevity.)**

So why all that attention paid to “forests”?



Angelic Forest

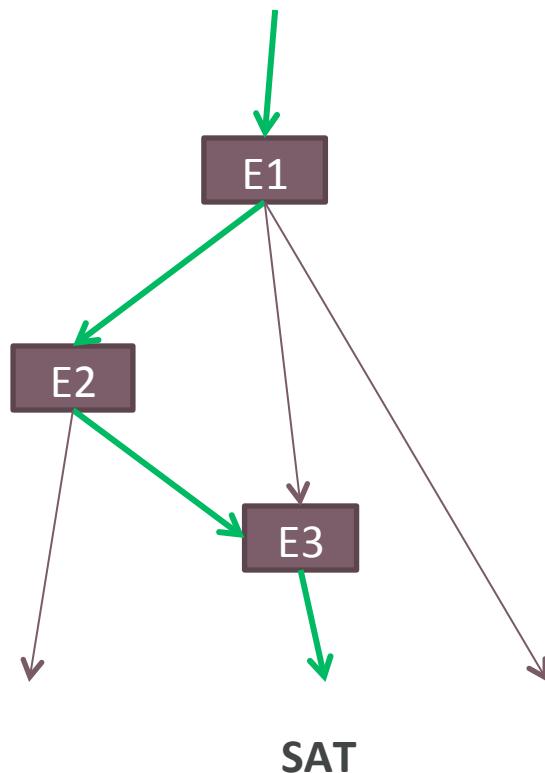
Program



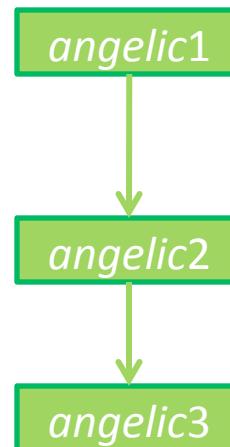
Angelic Paths

Angelic Forest

Program

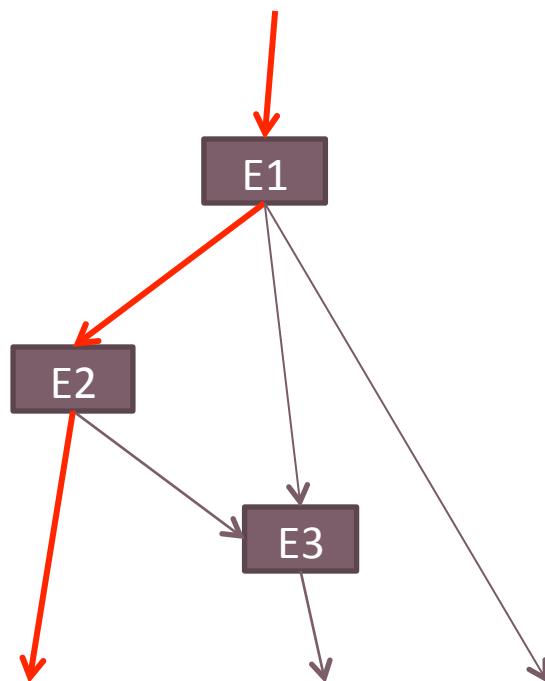


Angelic Paths

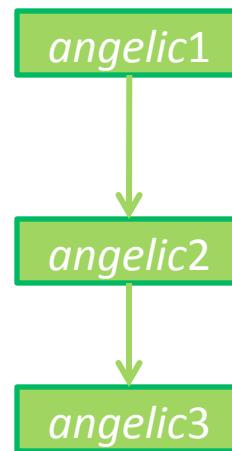


Angelic Forest

Program

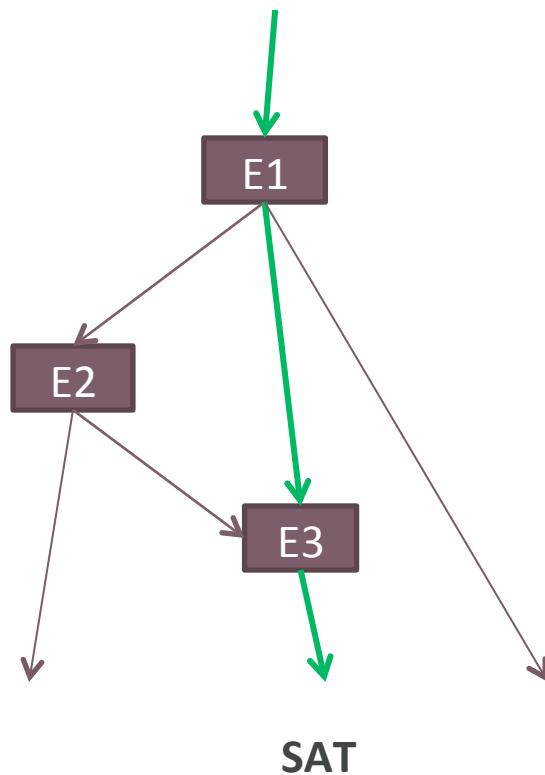


Angelic Paths

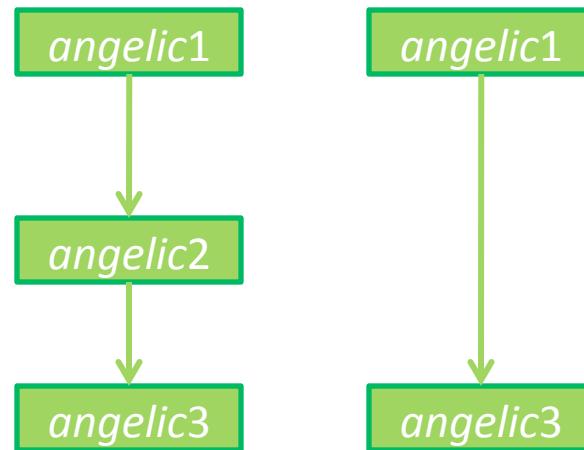


Angelic Forest

Program

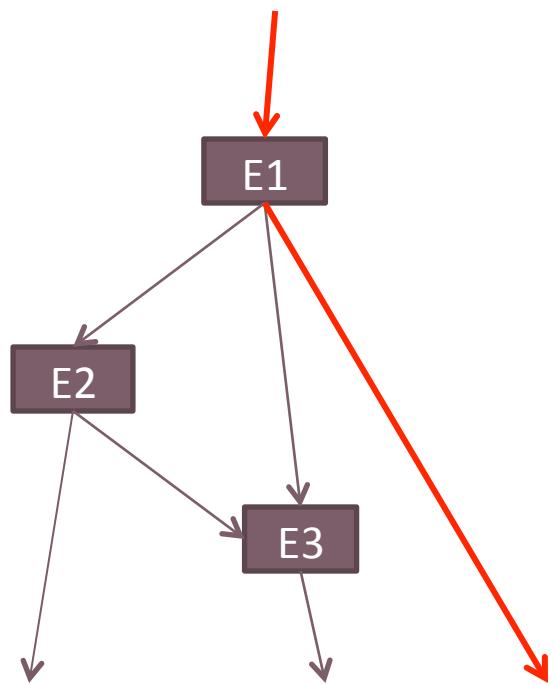


Angelic Paths

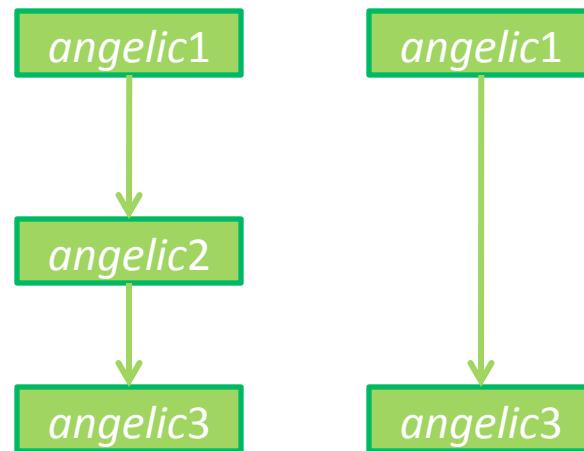


Angelic Forest

Program



Angelic Paths



UNSAT



Scalability

Tradeoffs and Challenges



Output quality

Expressive power



<https://www.flickr.com/photos/86530412@N02/7935377706>

: <https://pixabay.com/en/approved-control-quality-stamp-147677/>

<https://www.flickr.com/photos/cimmyt/5219256862>

Program	Description	LOC	Bug Type	Time (s)

Program	Description	LOC	Bug Type	Time (s)
gcd	example	22	infinite loop	153
nullhttpd	webserver	5575	heap buffer overflow (code)	578
zune	example	28	infinite loop	42
uniq	text processing	1146	segmentation fault	34
look-u	dictionary lookup	1169	segmentation fault	45
look-s	dictionary lookup	1363	infinite loop	55
units	metric conversion	1504	segmentation fault	109
deroff	document processing	2236	segmentation fault	131
indent	code processing	9906	infinite loop	546
flex	lexical analyzer generator	18774	segmentation fault	230
openldap	directory protocol	292598	non-overflow denial of service	665
ccrypt	encryption utility	7515	segmentation fault	330
lighttpd	webserver	51895	heap buffer overflow (vars)	394
atris	graphical game	21553	local stack buffer exploit	80
php	scripting language	764489	integer overflow	56
wu-ftpd	FTP server	67029	format string vulnerability	2256
leukocyte	computational biology	6718	segmentation fault	360
tiff	image processing	84067	segmentation fault	108
imagemagick	image processing	450516	wrong output	45 2160

Program	Description	LOC	Bug Type	Time (s)
gcd	example	22	infinite loop	153
nullhttpd	webserver	5575	heap buffer overflow (code)	578
zune	example	28	infinite loop	42
uniq	text processing	1146	segmentation fault	34
look-u	dictionary lookup	1169	segmentation fault	45
look-s	dictionary lookup	1363	infinite loop	55
units	metric conversion	1504	segmentation fault	109
deroff	document processing	2236	segmentation fault	131
indent	code processing	9906	infinite loop	546
flex	lexical analyzer generator	18774	segmentation fault	230
openldap	directory protocol	292598	non-overflow denial of service	665
ccrypt	encryption utility	7515	segmentation fault	330
lighttpd	webserver	51895	heap buffer overflow (vars)	394
atris	graphical game	21553	local stack buffer exploit	80
php	scripting language	764489	integer overflow	56
wu-ftpd	FTP server	67029	format string vulnerability	2256
leukocyte	computational biology	6718	segmentation fault	360
tiff	image processing	84067	segmentation fault	108
imagemagick	image processing	450516	wrong output	46 2160

Program	Description	LOC	Bug Type	Time (s)
gcd	example	22	infinite loop	153
nullhttpd	webserver	5575	heap buffer overflow (code)	578
zune	example	28	infinite loop	42
uniq	text processing	1146	segmentation fault	34
look-u	dictionary lookup	1169	segmentation fault	45
look-s	dictionary lookup	1363	infinite loop	55
units	metric conversion	1504	segmentation fault	109
deroff	document processing	2230	segmentation fault	131
indent	code processing	9906	infinite loop	546
flex	lexical analyzer generator	18774	segmentation fault	230
openldap	directory protocol	292598	non-overflow denial of service	665
ccrypt	encryption utility	7515	segmentation fault	330
lighttpd	webserver	51895	heap buffer overflow (vars)	394
atris	graphical game	21553	local stack buffer exploit	80
php	scripting language	764489	integer overflow	56
wu-ftpd	FTP server	67029	format string vulnerability	2256
leukocyte	computational biology	6718	segmentation fault	360
tiff	image processing	84067	segmentation fault	108
imagemagick	image processing	450510	wrong output	47 2160

Program	Description	LOC	Bug Type	Time (s)
gcd	example	22	infinite loop	153
nullhttpd	webserver	5575	heap buffer overflow (code)	578
zune	example	28	infinite loop	42
uniq	text processing	1146	segmentation fault	34
look-u	dictionary lookup	1169	segmentation fault	45
look-s	dictionary lookup	1363	infinite loop	55
units	metric conversion	1504	segmentation fault	109
deroff	document processing	2236	segmentation fault	131
indent	code processing	9906	infinite loop	546
flex	lexical analyzer generator	18774	segmentation fault	230
openldap	directory protocol	292598	non-overflow denial of service	665
ccrypt	encryption utility	7515	segmentation fault	330
lighttpd	webserver	51895	heap buffer overflow (vars)	394
atris	graphical game	21553	local stack buffer exploit	80
php	scripting language	764489	integer overflow	56
wu-ftpd	FTP server	67029	format string vulnerability	2256
leukocyte	computational biology	6718	segmentation fault	360
tiff	image processing	84067	segmentation fault	108
imagemagick	image processing	450516	wrong output	48 2160

**“IF I GAVE YOU THE LAST 100
BUGS FROM <MY PROJECT>,
HOW MANY COULD
GENPROG ACTUALLY FIX?”**

– MANY PEOPLE

Challenge: Indicative Bug Set



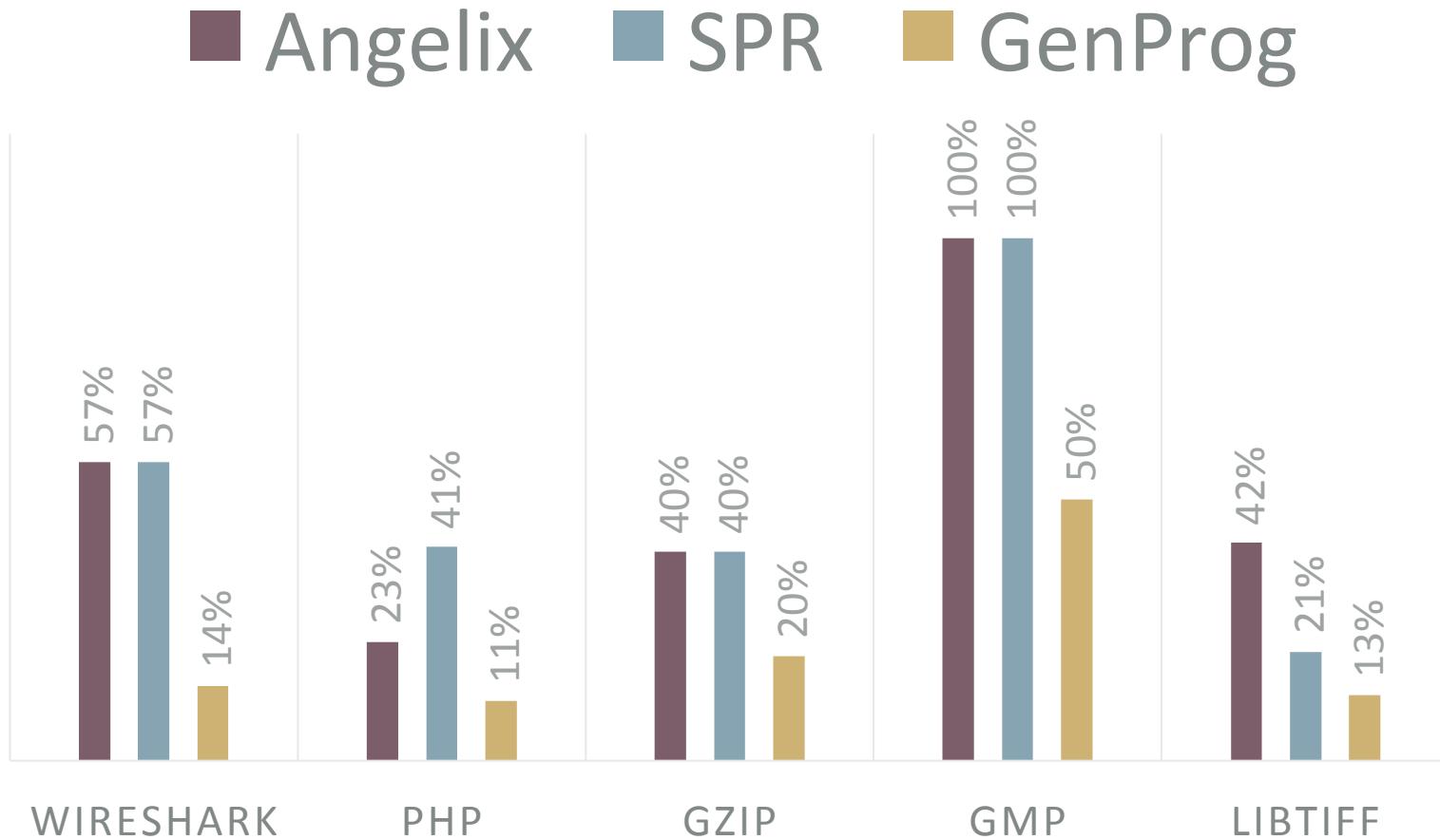
- Goal: a large set of **important, reproducible bugs** in **non-trivial** programs.
- Approach: use historical data (source control!) to approximate discovery and repair of bugs in the wild.

Success/Cost

Program	Defects Repaired	Cost per non-repair		Cost per repair	
		Hours	US\$	Hours	US\$
fbc	1/3	8.52	5.56	6.52	4.08
gmp	1/2	9.93	6.61	1.60	0.44
gzip	1/5	5.11	3.04	1.41	0.30
libtiff	17/24	7.81	5.04	1.05	0.04
lighttpd	5/9	10.79	7.25	1.34	0.25
php	28/44	13.00	8.80	1.84	0.62
python	1/11	13.00	8.80	1.22	0.16
wireshark	1/7	13.00	8.80	1.23	0.17
Total	55/105	11.22h		1.60h	

- \$403 for all 105 trials, leading to 55 repairs; \$7.32 per bug repaired.

Comparison (Repair-ability)



Heartbleed patch



NUS
National University
of Singapore

```
if (hbtype == TLS1_HB_REQUEST  
    && (payload + 18) < s->s3->rrec.length) {  
    ...  
} else if (hbtype == TLS1_HB_RESPONSE) {  
    ...  
}  
return 0;
```

Generated patch

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

Developer patch



Scalability

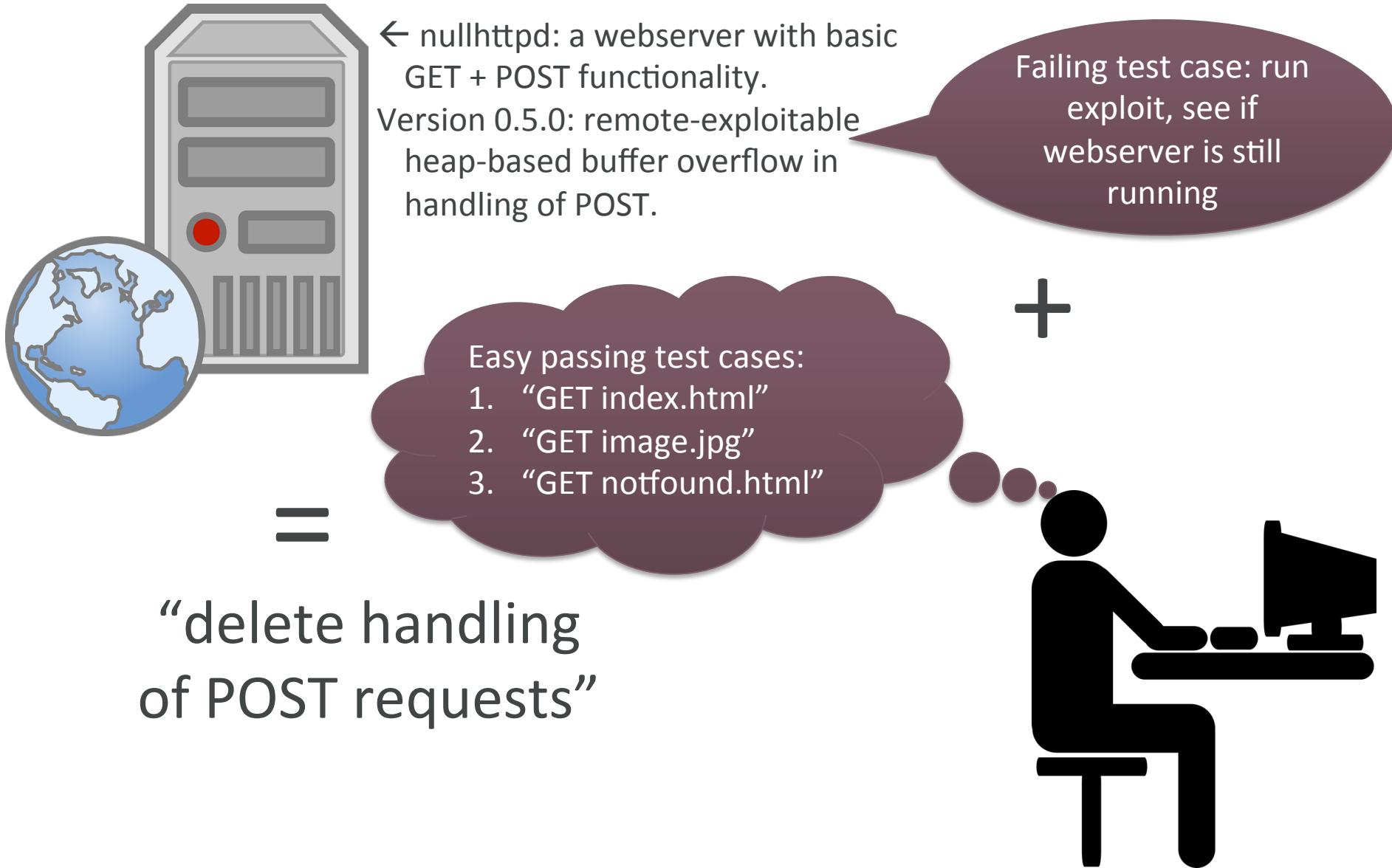
Tradeoffs and Challenges



Output quality

Expressive power





Flashback to 2008...

When we added a non-crashing test case for POST, proto-GenProg found a much better patch.

- When the test suite is your objective function, test suite quality matters.
 - ...how much is a trickier issue.
- But we're begging the question...



When we added a non-crashing test case for POST, proto-GenProg found a much better patch.

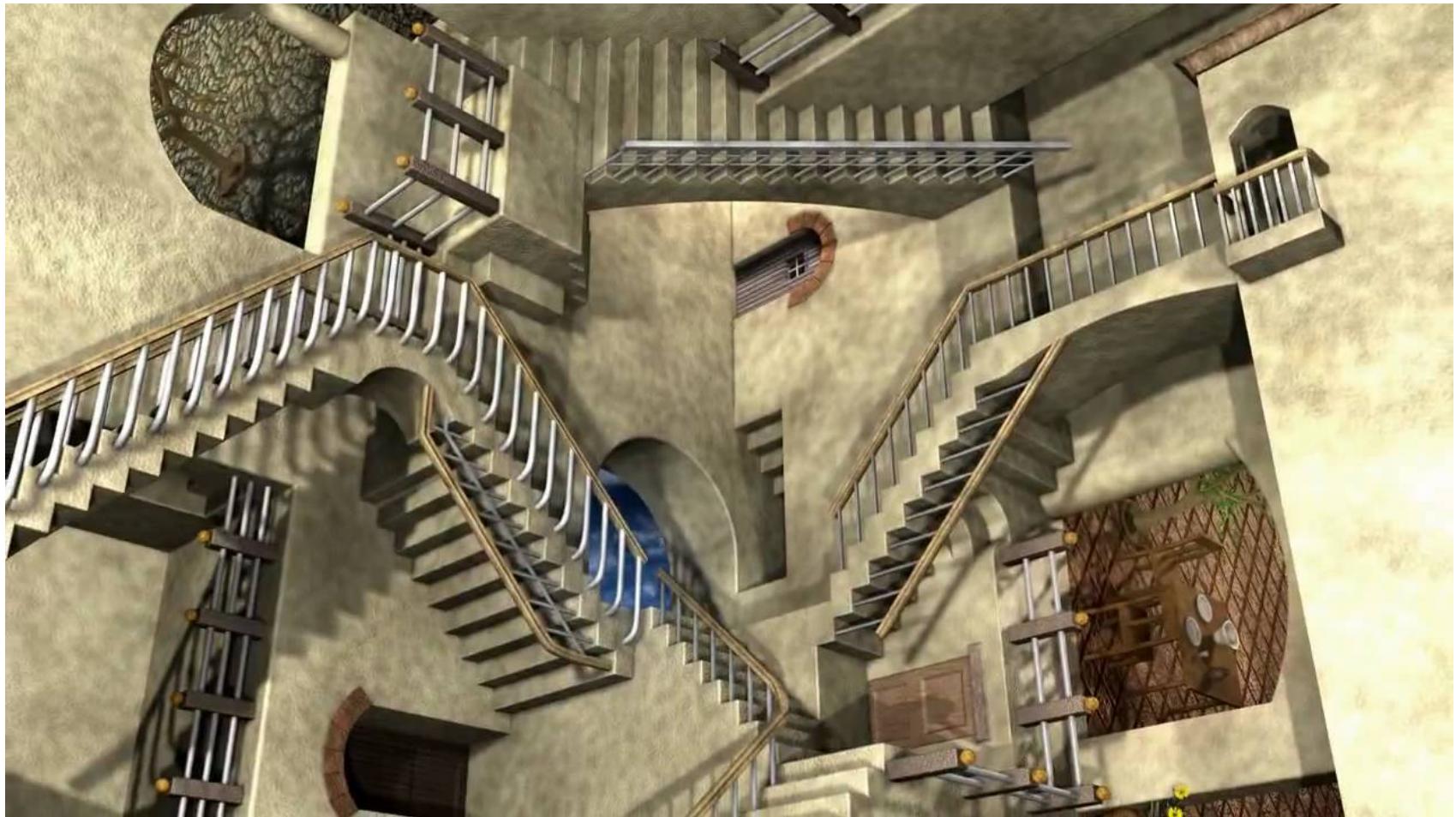
- When the test suite is your objective function, test suite quality matters.
 - ...how much is a trickier issue.
- But we're begging the question...

What is a high quality patch, anyway?

- Understandable?
 - Well, I had no problem understanding the POST-deleting patch...
 - (non-functional properties are important and being studied by others!)
- Doesn't delete?
 - But what about goto fail?
- Does the same thing the human did/would do?
 - But humans are often wrong! And how close does it have to be?
- Doesn't introduce new bugs?
 - How to tell?
- Addresses the cause, not the symptom...

Proposal: measure quality based on degree to which results *generalize*.

- In machine learning, techniques are trained and evaluated on disjoint datasets to assess overfitting.
- In program repair:
 - Tests used to build a repair are *training* tests
 - Tests used to assess correctness are *evaluation* tests



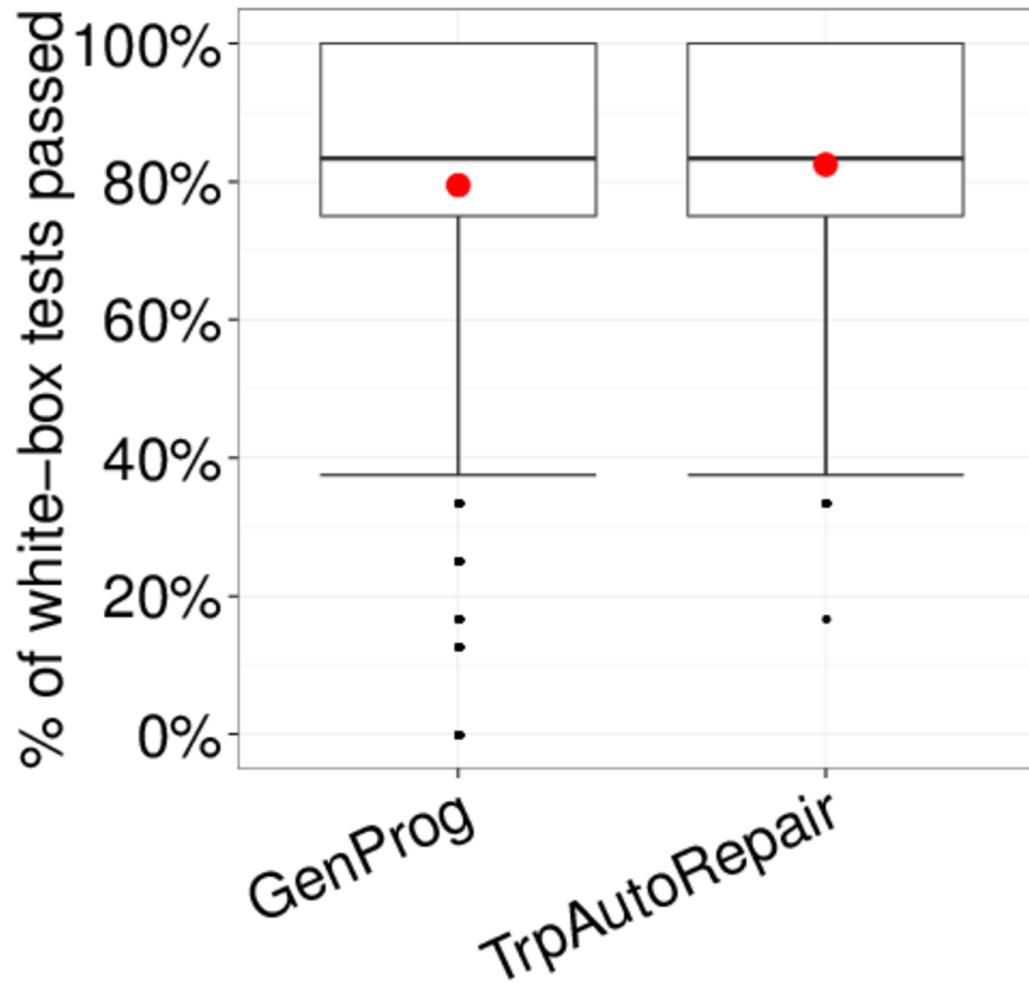
PROBLEM: THE DESIRED STUDY IS IMPOSSIBLE.

[Dataset + Tools]

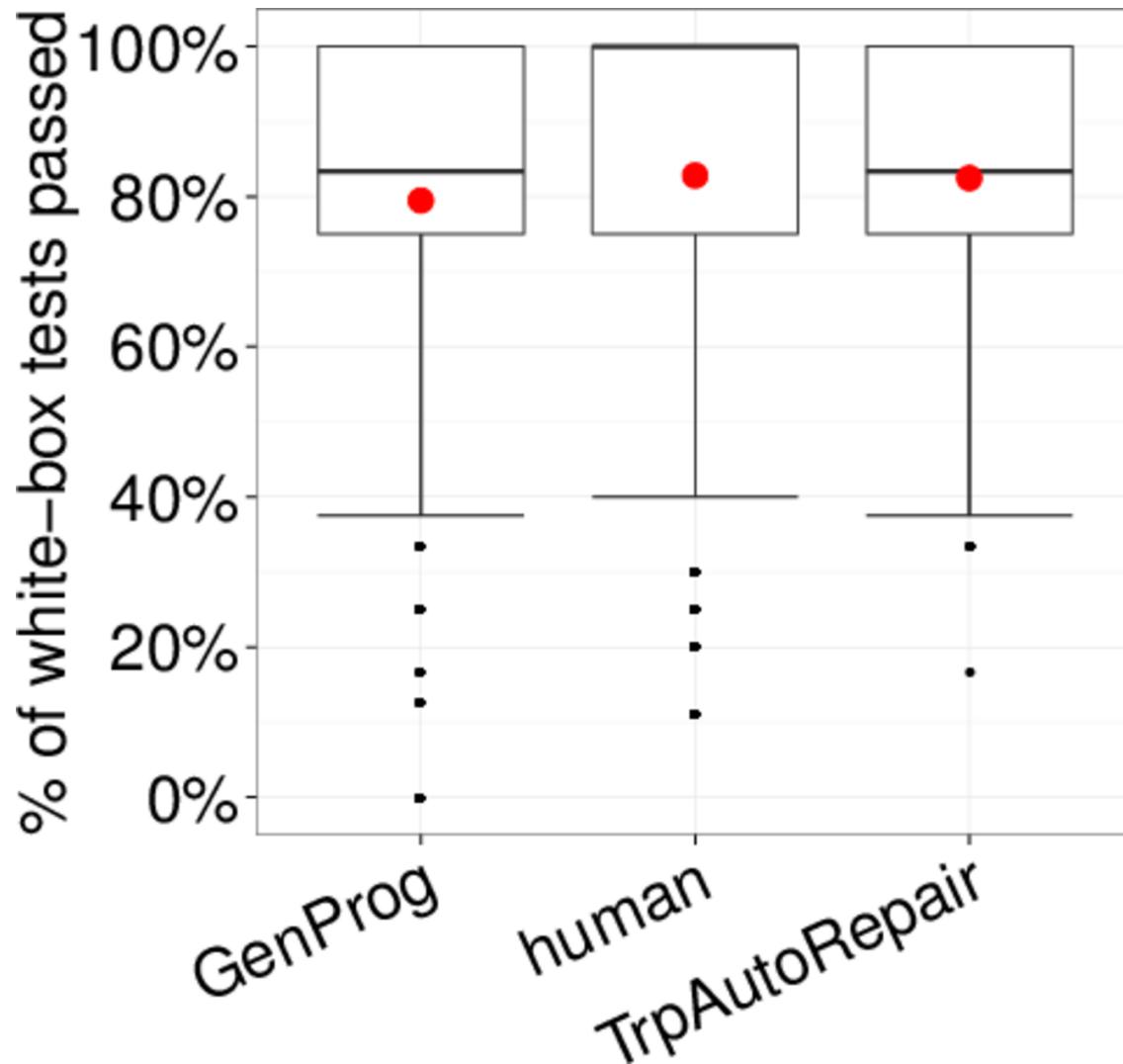
- Student homework submissions from six UC Davis Introduction to Programming assignments
- Two full-coverage test suites:
 - White-box suite generated by Klee from reference implementation.
 - Black-box suite written by course instructor.
 - *Feature: Assess patch quality as distinct from test suite quality.*
- Goal: Compare GenProg and TrpAutoRepair/RSRepair, G&V techniques with different search strategies.

Full dataset available at repairbenchmarks.cs.umass.edu

Both tools produced patches that overfit to the training set.



But: the tools do as well as the students!



Overfitting is not unique to heuristic techniques.

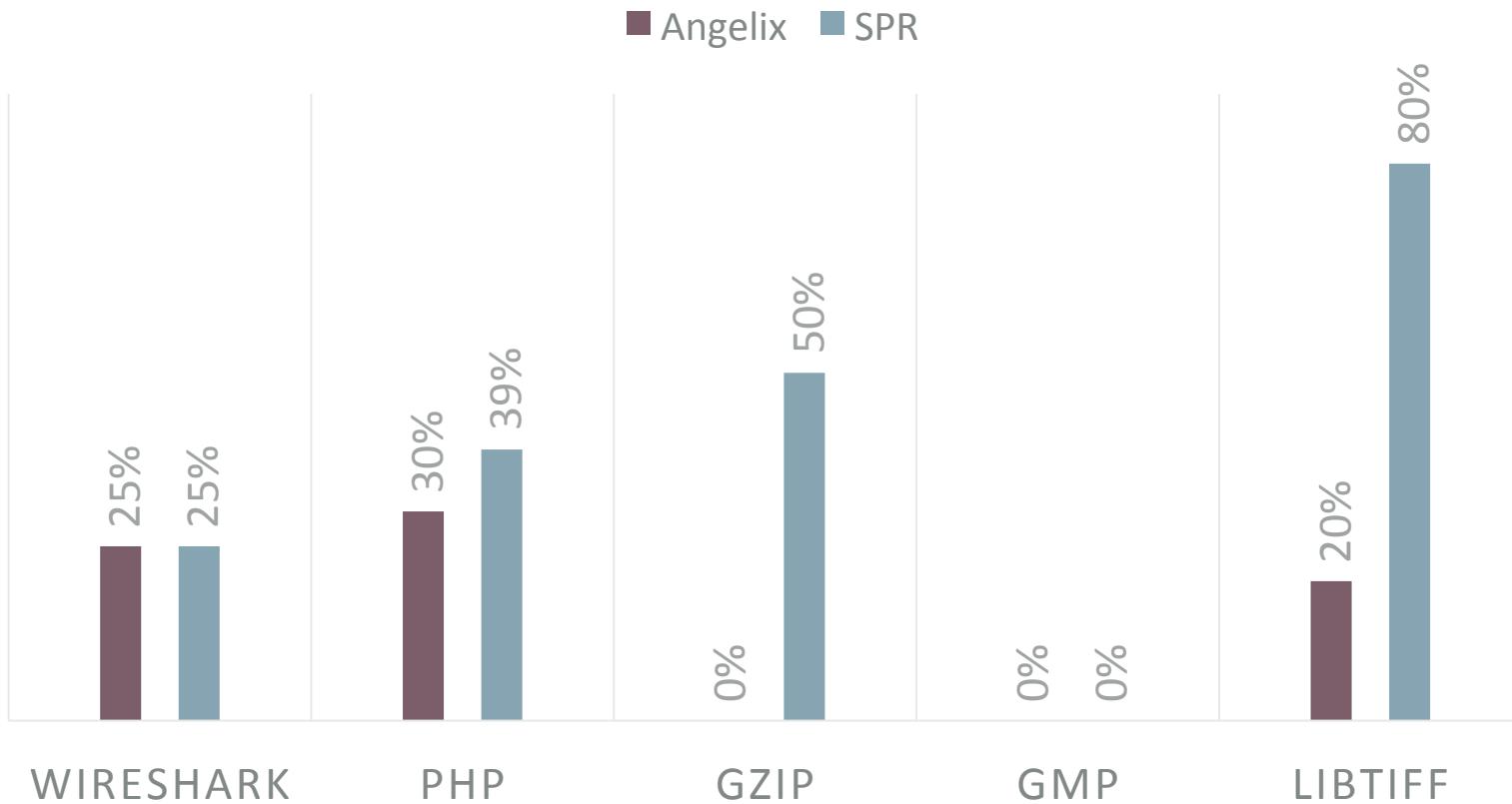
- Angelix: 120/233 of patches produced on a subset to IntroClass overfit.
- ~40% of SPR patches studied in Angelix paper delete functionality by generating tautological if conditions.

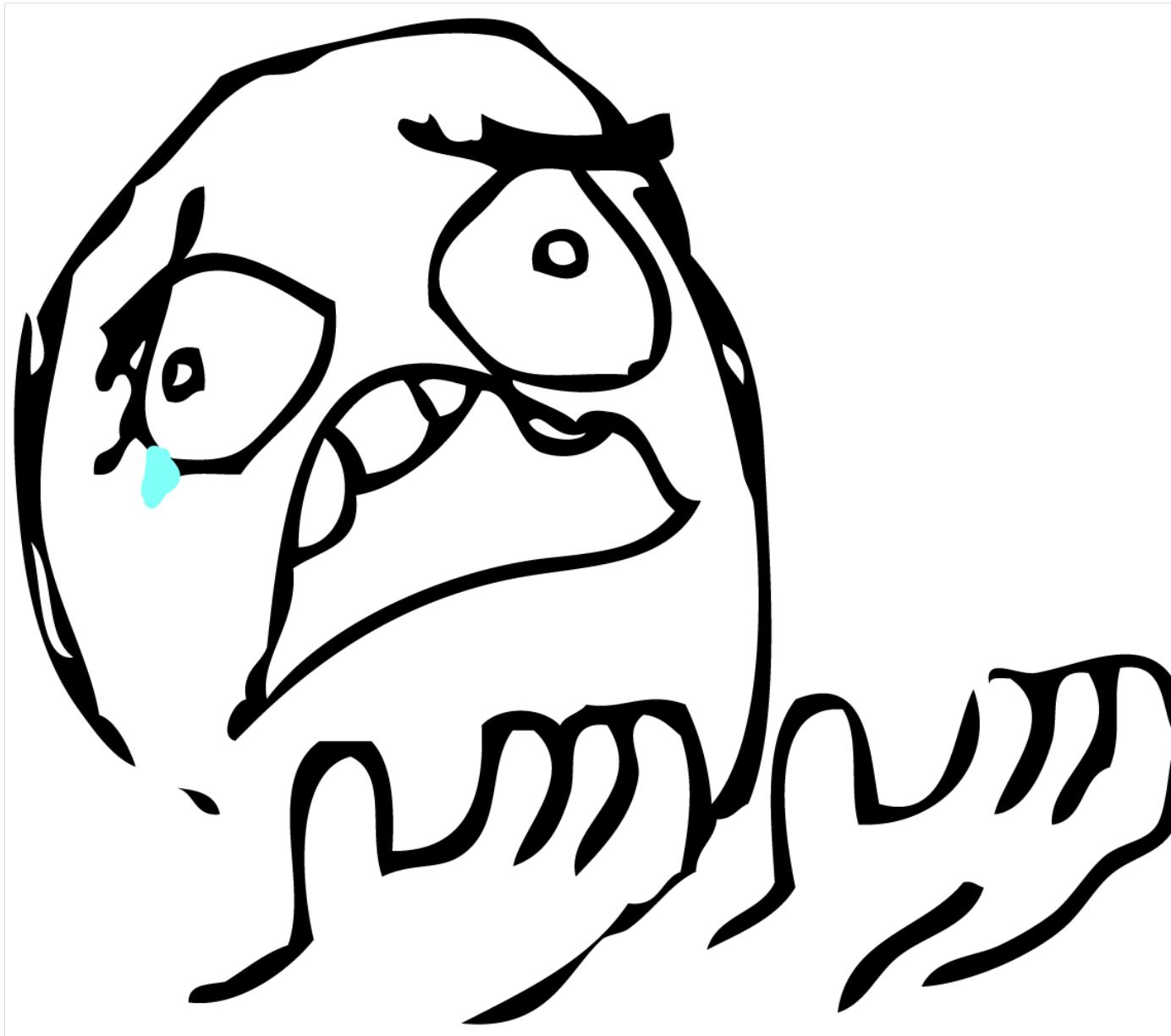
Overfitting is not unique to heuristic techniques.

- Angelix: 120/233 of patches produced on a subset of IntroClass overfit.
- ~40% of SMT papers do PhD students observe that much of the problem is in the synthesis of overly constrained tautological if conditions.

Quality Comparison with SPR

FUNCTIONALITY-DELETING REPAIRS





Context matters!

**OPTION 1: UNDERSTAND AND REASON
ABOUT THE CIRCUMSTANCES UNDER
WHICH PERFECTION IS NOT REQUIRED.**



← Scenario: Long-running servers +
IDS + generate repairs for
detected anomalies.

Workloads: a day of unfiltered
requests to the LIVA CC
webserver

**THIS PATCH
WAS BAD**

**NullHttpd
(logo-less)**

php

LIGHTTPD
fly light.

2012 flashback...

Even a functionality-reducing repair had little practical impact.

Program	Post-patch requests lost	Fuzz Tests Failed	
		General	Exploit
nullhttpd	$0.00 \% \pm 0.25\%$	$0 \rightarrow 0$	$10 \rightarrow 0$
lighttpd	$0.03\% \pm 1.53\%$	$1410 \rightarrow 1410$	$9 \rightarrow 0$
php-BAD	$0.02\% \pm 0.02\%$	$3 \rightarrow 3$	$5 \rightarrow 0$

How?

**OPTION 2: DEVELOP TECHNIQUES THAT
ARE MORE LIKELY TO GENERALIZE.**

Challenge your assumptions!



**EXAMPLE
ASSUMPTION:
bug-fixing
patches are like
kittens: smaller
is better!**

What if...

- Instead of trying to make small changes, we replaced buggy regions with code that correctly captures the overall desired logic?
- Principle: using human-written code to fix code at a higher granularity level leads to better quality repairs.

SEARCHREPAIR:

HIGH-QUALITY

AUTOMATED BUG

REPAIR USING

SEMANTIC SEARCH

Semantic code search looks for code based on what it should *do*.

- Keyword: “C median three numbers”
- Semantic:

Input	Expected
2,6,8	6
2,8,6	6
6,2,8	6
6,8,2	6
8,6,2	6
9,9,9	9

...Generate and validate + Semantic reasoning!

SearchRepair patches were of much higher quality than those produced by previous techniques.

Technique	Held out tests passed
SearchRepair	97.2%
GenProg	68.7%
TRPAP	72.1%
AE	64.2%





Scalability

The Three Major Challenges



Output quality

Expressive power





Young Claire



Line 7

inhibit = 1, up_sep

bias = α , PC= $\alpha > 110$



Line 8

down_sep = 110

bias = α , PC= $\alpha \leq 110$





COLLABORATORS MAKE THE WORLD GO ROUND

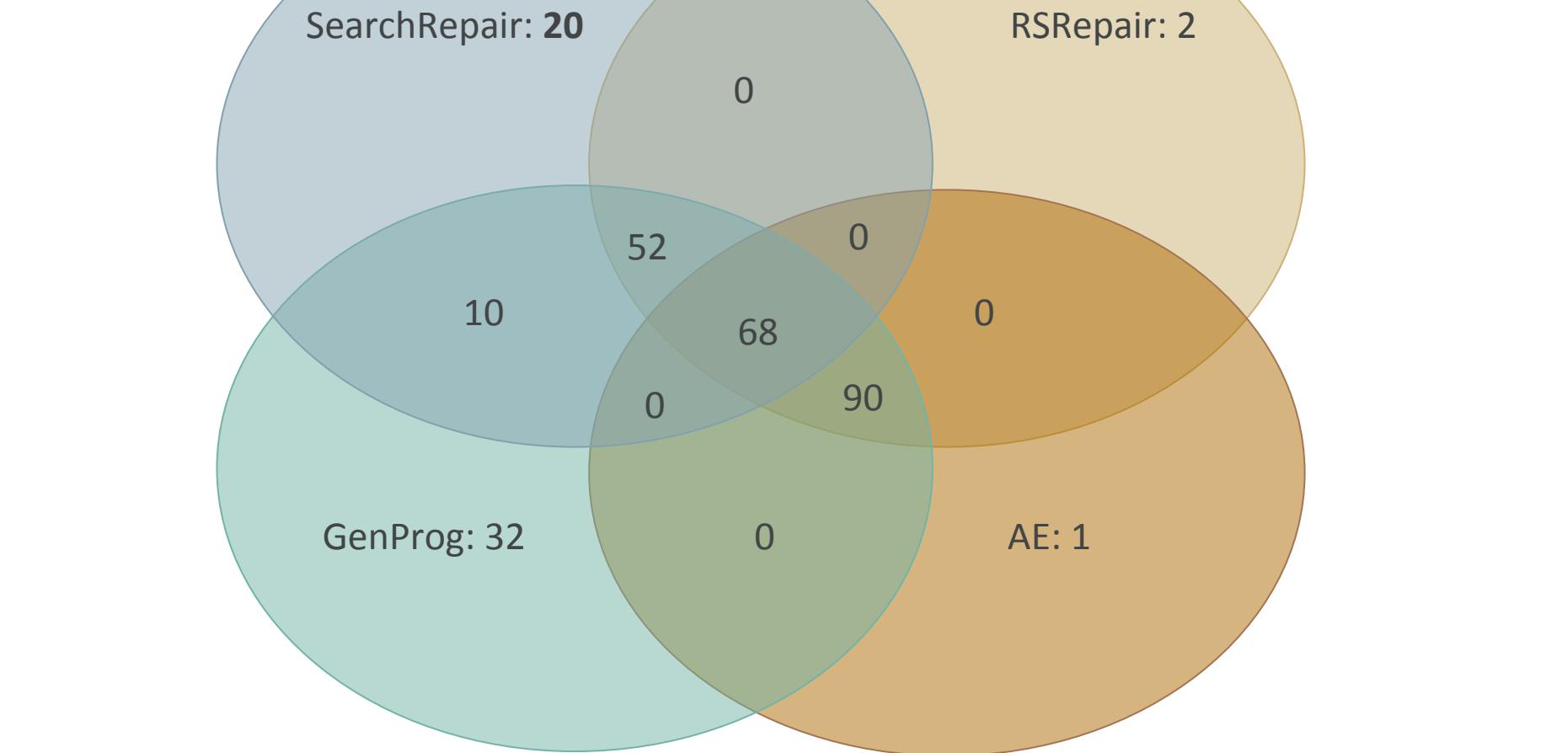


SearchRepair total: 150

RSRepair total: 247

SearchRepair: 20

RSRepair: 2



GenProg total: 287

AE total: 159