# Concolic Program Repair

Detecting and discarding over-fitting patches via systematic co-exploration of the patch space and input space

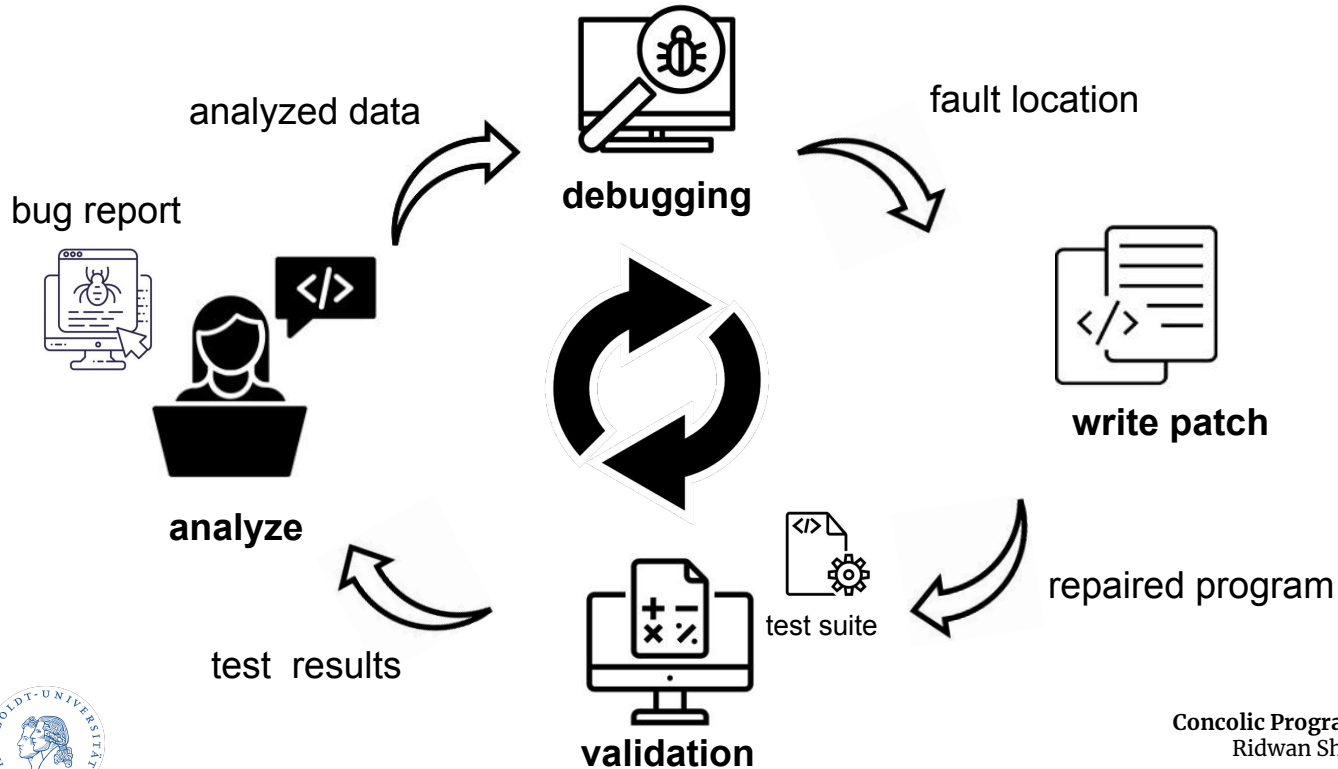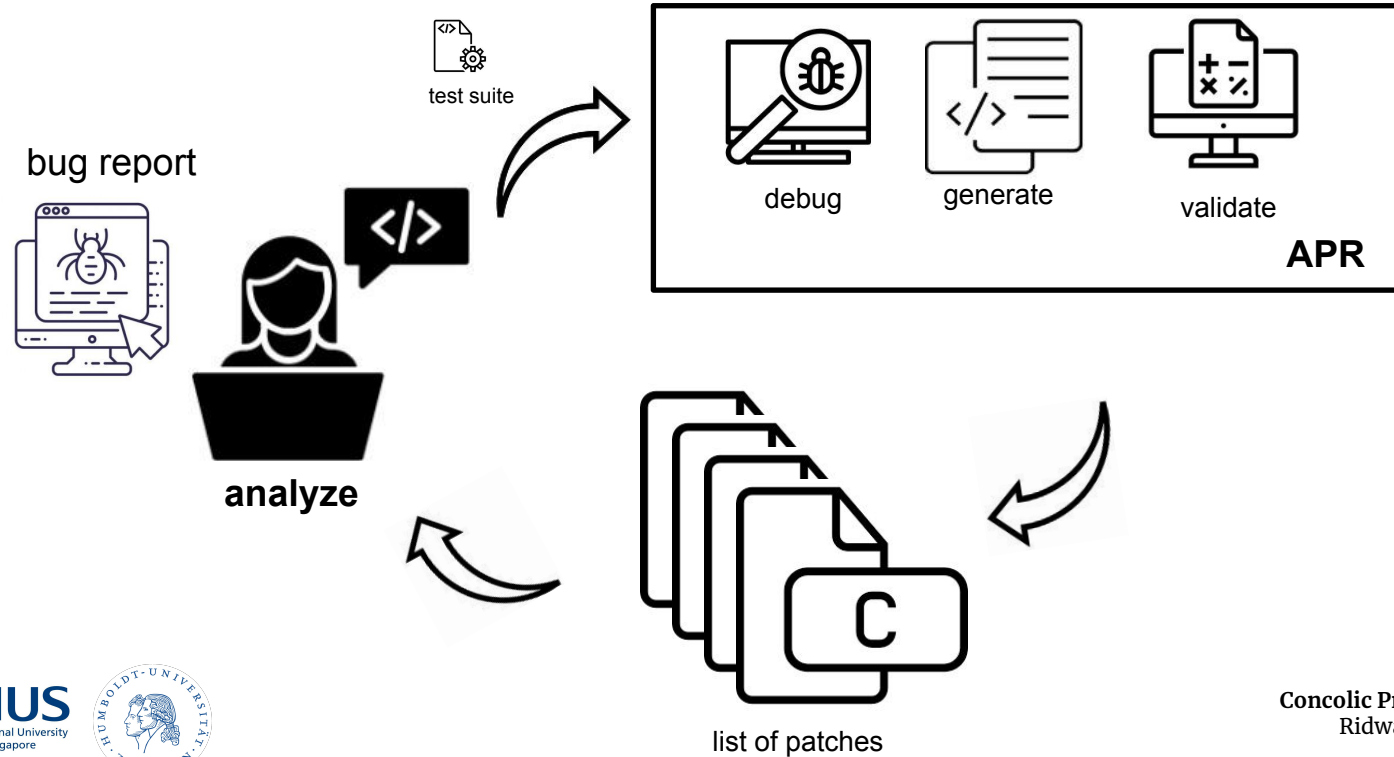**Ridwan Shariffdeen**   **Yannic Noller**   **Lars Grunske**   **Abhik Roychoudhury**

# Typical Repair Workflow



analyzed data

fault location

bug report

**debugging**

**write patch**

**analyze**

repaired program

test suite

test results

**validation**

# Automated Program Repair



bug report

test suite

debug    generate    validate

**APR**

**analyze**

list of patches

# Trust Enhancement Issues in Program Repair

Yannic Noller, **Ridwan Shariffdeen**, Xiang Gao, Abhik Roychoudhury

**IEEE/ACM 44th International Conference on Software Engineering (ICSE) 2022**

# Developer Survey – Demographics

**103** software practitioners

89% with **2+ years experience**

75% **Software Developers**

35 questions on trust in APR

# Insights from Developer Survey

**I**    **Additional test-cases** improves trustworthiness of generated pa

**II**    Developers are willing to **provide specification** to the repair pro

**III**    Developers will only allow a maximum of **1-hour timeout**

**IV**    Developers will only review up to **maximum of 5 patches**

https://zenodo.org/record/6303481

**Concolic Program Repair**
Ridwan Shariffdeen

6

# Concolic Program Repair

**Ridwan Shariffdeen,** Yannic Noller, Lars Grunske, Abhik Roychoudhury

**42nd ACM SIGPLAN Conference on Programming Language Design and Implementation (PLDI) 2021**

# Key Idea: Gradual Correctness

Detecting and discarding over-fitting patches via systematic co-exploration of the patch space and input space

# Our Solution

**semantic approach** incl. **program synthesis**
- ➡ avoids **non-compilable** patches
- ➡ provides **symbolic reasoning** capabilities

**co-exploration** of the **input** space and **patch** space
- ➡ prune **over-fitting** patches
- ➡ enables **gradual improvement**

**user-provided specification**
- ➡ to **reason** about **additional inputs**
- ➡ **key aspect** to handle absence of test cases

**Input Space**

initial test case

P3
P2
P1
P4

explored path
(input partition)

**infeasbility** checks
in both directions

**Patch Space**

refined patch set

plausible
patches

correct
patch (set)

represented with
**abstract patches**

**Concolic Program Repair**

NUS
National University
of Singapore

HUMBOLDT-UNIVERSITÄT
ZU BERLIN

# Patch Representation

## .. **concrete** patches

```
x > 0
x > 1
x > 2
…
```

```
x + 1 > y
x - 1 > y
x + 2 > y
…
```

## .. **abstract** patches

```
x > a, a ∈ [0, 10]
```

```
x + a > y, a ∈ [-10, 10]
```

Our notion of an **abstract patch** represents a **patch template** with **parameters.**

➡ **generate** and **maintain** smaller amount of patch candidates

➡ allows refinement instead of just discarding

➡ **subsumes** concrete patches

# Abstract Patches

$$(\boldsymbol{\theta_\rho}, \boldsymbol{T_\rho}, \boldsymbol{\psi_\rho})$$

$X_\rho$ is the set of **program variables**
$X \subseteq X_\rho$ is the set of **input variables**
$A$ is the set of **template parameters**

- $\theta_\rho(X_\rho, A)$ denotes the **repaired** (boolean or integer) **expression**

- $T_\rho(A)$ represents the **conjunction of constraints** $\tau_\rho(a_i)$ on the **parameters** $a_i \in A$ included in $\theta_\rho$: $T_\rho(A) = \wedge_{a_i \in A} \tau_\rho(a_i)$

- $\psi_\rho(X, A)$ is the **patch formula induced by inserting** the expression $\theta_\rho$ into the buggy program

1. patch is a **condition**

```
        uint32 rnrows = roundup(nrows,ve
        if (CONDITION) return 0;
        /* potential divide-by-zero err
```

```
if (ρ)
    return 0;
```

$\theta_\rho := x > a$
$T_\rho = \tau_\rho(a) := (a \geq -10)$
$\psi_\rho := x > a$

2. patch is a **right hand-side of an assignment**

```
...
y = ρ;
...
```

$\theta_\rho := x - a$
$T_\rho = \tau_\rho(a) := (a \geq -10)$
$\psi_\rho := (y = x - a)$

# New kind of Infeasibility

## .. in the **input space**

**Path Reduction:**
For every generated input, we check that there is one patch that can exercise the corresponding path. Otherwise, the path will not be explored.

For example:

$$\phi := x > 3 \; \wedge \; y > 5 \; \wedge \; \rho$$
$$\rho := (x = 0 \; \vee \; y = 0)$$

## .. in the **patch space**

**Patch Reduction:**
If a patch allows inputs to exercise a path that violates the specification, we identify this as a patch that overfits the valid set of values and attempt to refine it.

parameters    inputs

$$\forall a_1, a_2, .., a_n \; \forall x_1, x_2, .., x_m :$$
$$\phi(X) \wedge \psi_\rho(X, A) \wedge T_\rho(A) \implies \sigma(X)$$

path constraint  patch constraint  parameter constraint  specification

# Our Implementation

# Illustrative Example



## Input

Buggy Program

Failing test case(s)

Fix Locations

User Specification

CVE-2016-3623:
Divide by Zero in LibTIFF v4.0.6

e.g., exploit as
TIFF picture

source location, (fix template),
synthesis components

formula about correct
behavior in SMT format

```
........
static int
cvtRaster(TIFF* tif, uint32* raster, uint32 width...)
{
    uint32 y;
    tstrip_t strip = 0;
    tsize_t cc, acc;
    unsigned char* buf;
    uint32 rwidth = roundup(width, horizSubSampling);
    uint32 rheight = roundup(height, vertSubSampling);
    uint32 nrows = (rowsperstrip > rheight ?
        rheight : rowsperstrip);
    uint32 rnrows = roundup(nrows,vertSubSampling);
    if (CONDITION) return 0;
    /* potential divide-by-zero error */
    cc = rnrows*rwidth + 2 * ((rnrows*rwidth)
        / (horizSubSampling*vertSubSampling));
    ........
}
```

observation

(assert (= false (= observation 0)))

# Illustrative Example

```
uint32 nrows = roundup(nrows,v
if (CONDITION) return 0;
/* potential divide-by-zero err
```

x ≜ horizSubSampling
y ≜ vertSubSampling
C ≜ CONDITION

## Input Space

**I**

P3
P2
P1
P4

Initial test input
x=7, y=0

**II**

P3
P2
P1
P4

P1: x > 3 ∧ y ≤ 5 ∧ ¬C

## Patch Space

69

plausible patches

correct patch (set)

46

## Patch Details

| ID | Patch Template | Parameter Constraint | # Conc. Patches |
|---|---|---|---|
| 1 | x >= a | a ≥ -10 ∧ a ≤ 7 | 18 |
| 2 | y < b | b ≥ 1 ∧ b ≤ 10 | 10 |
| 3 | x == a ‖ y == b | (a=7 ∧ b ≥ -10 ∧ b ≤ 10) ∨ (b=0 ∧ a ≥ -10 ∧ a ≤ 10) | 41 |

| ID | Patch Template | Parameter Constraint | # Conc. Patches |
|---|---|---|---|
| 1 | x >= a | a ≥ -10 ∧ a ≤ 4 | 15 |
| 2 | y < b | b ≥ 1 ∧ b ≤ 10 | 10 |
| 3 | x == a ‖ y == b | b=0 ∧ a ≥ -10 ∧ a ≤ 10 | 21 |

NUS National University of Singapore

HUMBOLDT-UNIVERSITÄT ZU BERLIN

# Illustrative Example



**Input Space**

**Patch Space**

**Patch Details**

I — Initial test input x=7, y=0 — P3, P2, P1, P4 — correct patch, plausible patches — 69

| ID | Patch Template | Parameter Constraint | # Conc. Patches |
|---|---|---|---|
| 1 | x >= a | a ≥ -10 ∧ a ≤ 7 | 18 |
| 2 | y < b | b ≥ 1 ∧ b ≤ 10 | 10 |
| 3 | x == a \|\| y == b | (a=7 ∧ b ≥ -10 ∧ b ≤ 10) ∨ (b=0 ∧ a ≥ -10 ∧ a ≤ 10) | 41 |

II — P3, P2, P1, P4 — P1: x > 3 ∧ y ≤ 5 ∧ ¬C — 46

| ID | Patch Template | Parameter Constraint | # Conc. Patches |
|---|---|---|---|
| 1 | x >= a | a ≥ -10 ∧ a ≤ 4 | 15 |
| 2 | y < b | b ≥ 1 ∧ b ≤ 10 | 10 |
| 3 | x == a \|\| y == b | b=0 ∧ a ≥ -10 ∧ a ≤ 10 | 21 |

III — P3, P2, P1, P4 — P2: x ≤ 3 ∧ y > 5 ∧ ¬C — 12

| ID | Patch Template | Parameter Constraint | # Conc. Patches |
|---|---|---|---|
| 1 | x >= a | a ≥ -10 ∧ a ≤ 0 | 11 |
| 2 | y < b | False | 0 |
| 3 | x == a \|\| y == b | a = 0 ∧ b = 0 | 1 |

IV — P3, P2, P1, P4 — P3: x ≤ 3 ∧ y ≤ 5 ∧ ¬C — 1

| ID | Patch Template | Parameter Constraint | # Conc. Patches |
|---|---|---|---|
| 1 | x >= a | False | 0 |
| 3 | x == a \|\| y == b | a = 0 ∧ b = 0 | 1 |

V — P3, P2, P1, P4 — P4: x > 3 ∧ y > 5 ∧ C — 1

| ID | Patch Template | Parameter Constraint | # Conc. Patches |
|---|---|---|---|
| 3 | x == a \|\| y == b | a = 0 ∧ b = 0 | 1 |

# Evaluation Setup

**Techniques**

CEGIS
ExtractFix
Angelix
Prophet

**Benchmarks**

ExtractFix

ManyBugs

SV-COMP

**Focus Areas**

Vulnerability Repair

Test-based Repair

Fixing Assertions

# Evaluation Insights

| Program | #Vul | Prophet | Angelix | ExtractFix | CPR |
|---------|------|---------|---------|------------|-----|
| LibTIFF | 11 | 1 | 0 | 6 | 11 |
| Binutils | 2 | - | - | 1 | 2 |
| LibXML2 | 5 | 0 | 0 | 2 | 5 |
| LibJPEG | 4 | 1 | - | 2 | 4 |
| FFmpeg | 2 | - | - | 2 | - |
| Jasper | 2 | 0 | 0 | 1 | 1 |
| Coreutils | 4 | 0 | - | 2 | 4 |
| **Total** | 30 | 2 | 0 | 16 | **27** |

Number of correct patches generated
for ExtractFix benchmark in 1h timeout

Up **74%** Patch Space Reduction

CPR **can gradually refine** the patch space via concolic exploration

CPR can be used for **test-guided general-purpose repair** and **security repair**

CPR is **more effective** than CEGIS wrt input and space exploration

NUS
National University
of Singapore

HUMBOLDT-UNIVERSITÄT ZU BERLIN

# Artifact



https://cpr-tool.github.io

http://doi.org/10.5281/zenodo.4668317