

# Hybrid Differential Software Testing

Disputation

**Yannic Noller**

Humboldt-Universität zu Berlin  
yannic.noller@acm.org

# Agenda

Problem

1. Problem + Motivation: Differential Testing

Contribution

2. Contribution Summary

Background

3. Foundations: Fuzzing + Symbolic execution

4. Related Work: Differential Testing

Solutions

5. Solutions:

Differential Fuzzing

Differential Dynamic Symbolic Execution

HyDiff (Hybrid Differential Software Testing)

Validation

6. Validation

Summary

7. Conclusion

# Software Engineering

”systematic application of scientific and technological knowledge, methods, and experience to the design, implementation, **testing**, and documentation of software”

[IEEE2017]

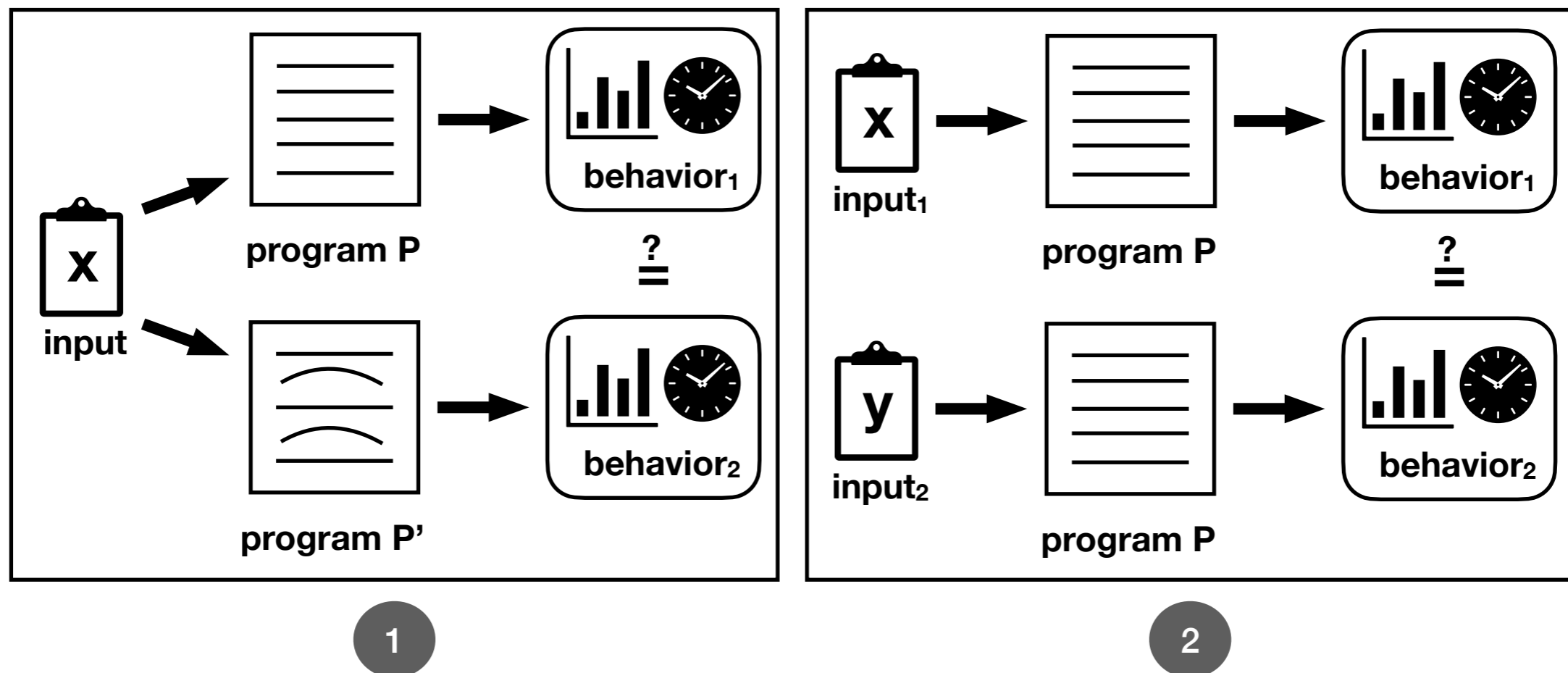
**Software Quality  
Assurance**



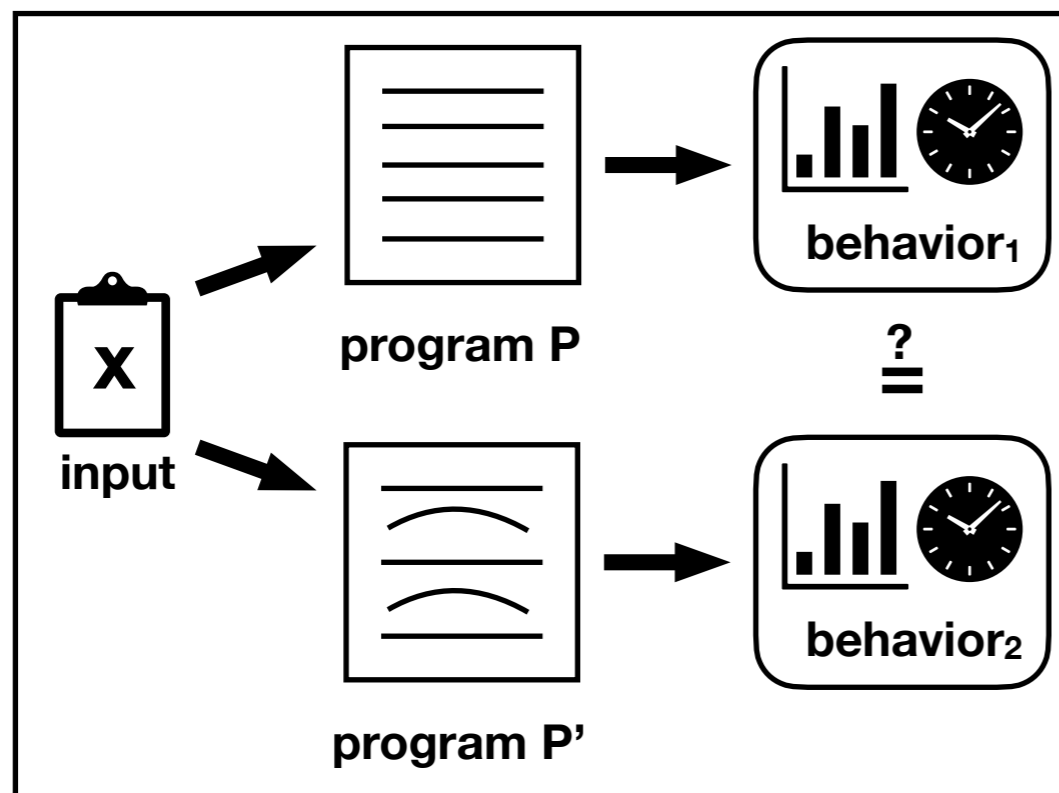
**Software  
Testing**



# Differential Software Testing



# Differential Software Testing



1

between **two program versions**  
for the **same input**  
↳ *software maintenance*

**Regression Analysis**

# Regression Analysis

```

1  int foo (int x) {
2      int y;
3      if (x < 0) {
4-         y = -x;
4+         y = x * x;
5      } else {
6         y = 2 * x;
7      }
8+     y = y + 1;
9     if (y > 1) {
10        return 0;
11    } else {
12        if (y == 1)
13            assert(false);
14    }
15    return 1;
16 }

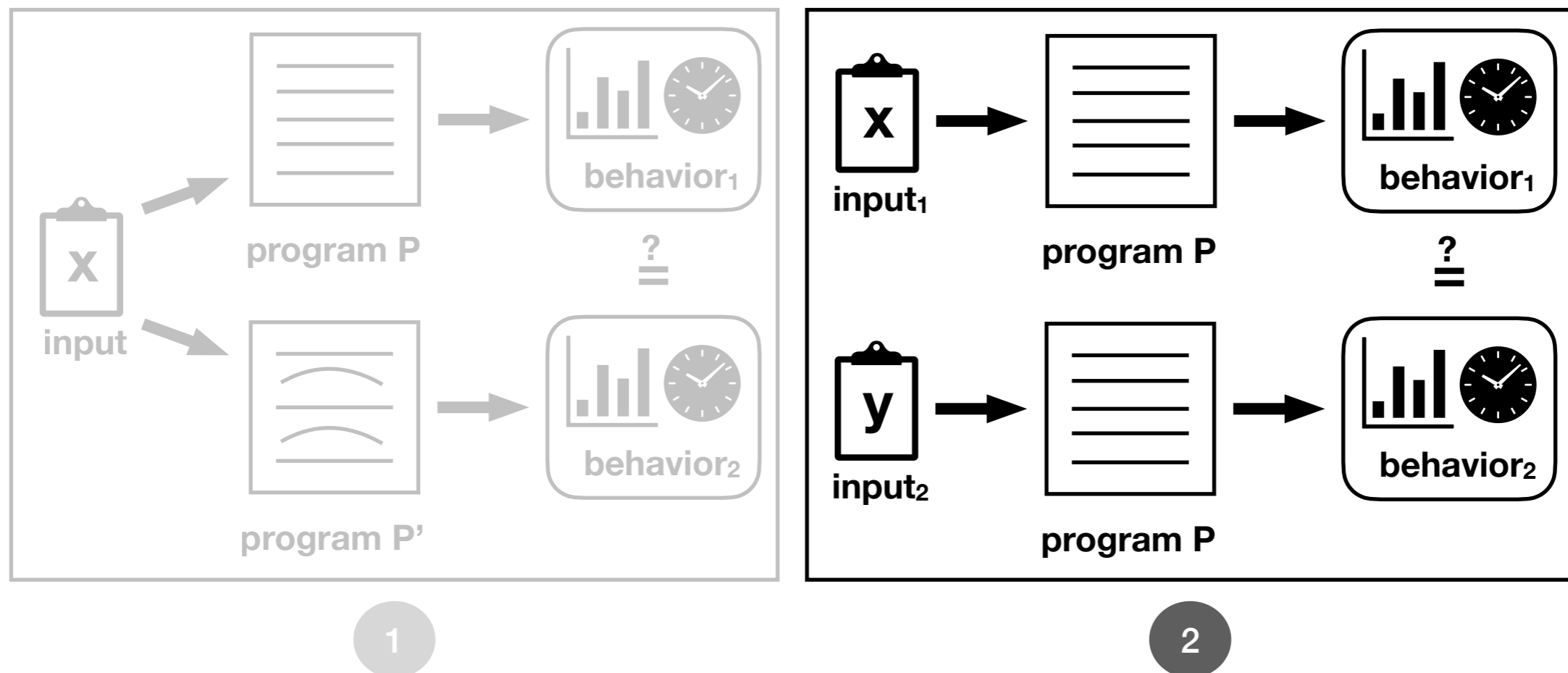
```

*Fixed* assertion error  
for  $x=-1$  (returns 0).

Are there unintended  
**behavioral** differences  
between the two  
versions?

*introduced new*  
assertion error  
for  $x=0$   
(previously returned 1)  
→ **Regression Bug**

# Differential Software Testing





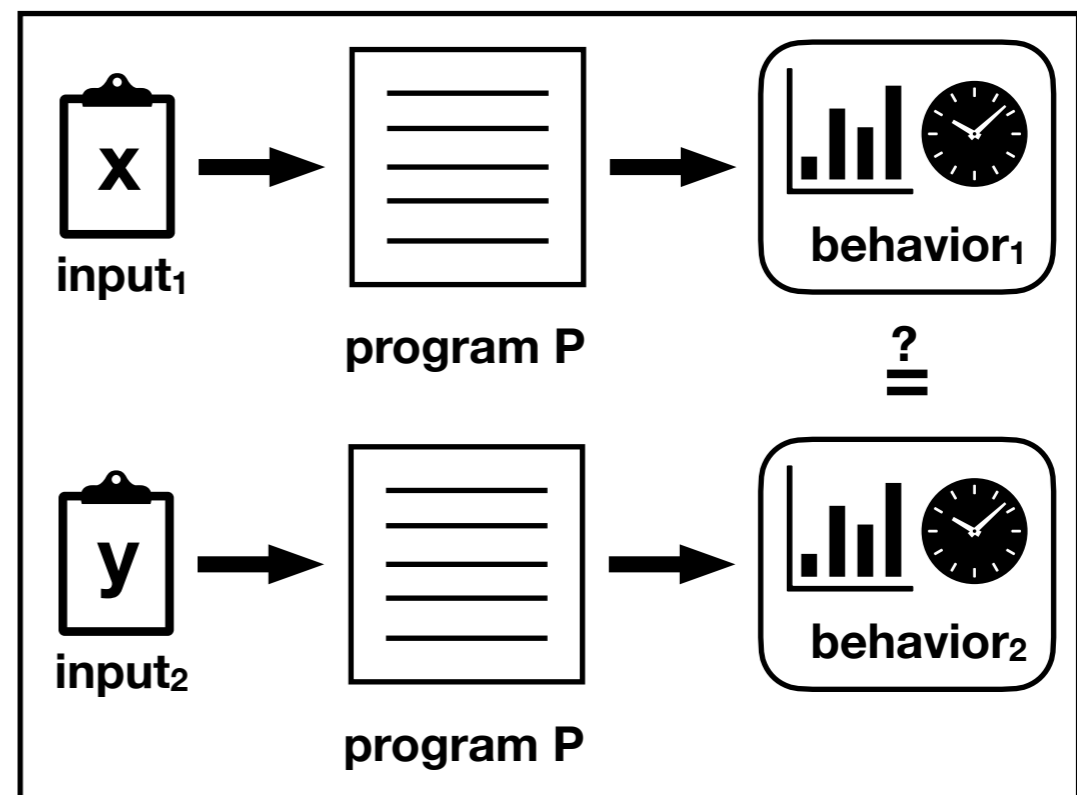
# Differential Software Testing

for the **same program** with  
**two different inputs**  
↳ *security, reliability*

Worst-Case Complexity  
Analysis

Side-Channel Analysis

Robustness Analysis of  
Neural Networks



2

# Worst-Case Complexity Analysis

**Goal:** discover vulnerabilities related to algorithmic complexity

```
0 public void sort (int[] a) {  
1     int N = a.length;  
2     for (int i = 1; i < N; i++) {  
3         int j = i - 1;  
4         int x = a[i];  
5         while ((j >= 0) && (a[j] > x)) {  
6             a[j + 1] = a[j];  
7             j--;  
8         }  
9         a[j + 1] = x;  
10    }  
11 }
```

## Insertion Sort

➔ find worst-case input:  
automated + fast + concrete

- worst-case complexity:  
 $O(n^2)$
- e.g.  $a=[8, 7, 6]$  ( $n=3$ )

# Side-Channel Analysis

- leakage of **secret** information
- **software** side-channels
- **observables:**
  - execution time,
  - memory consumption,
  - response size,
  - ...

# Example: Side-Channel Vulnerability

```
0  boolean pwcheck_unsafe (byte[] pub, byte[] sec) {  
1      if (pub.length != sec.length) {  
2          return false;  
3      }  
4      for (int i = 0; i < pub.length; i++) {  
5          if (pub[i] != sec[i]) {  
6              return false;  
7          }  
8      }  
9      return true;  
10 }
```

## Unsafe Password Checking

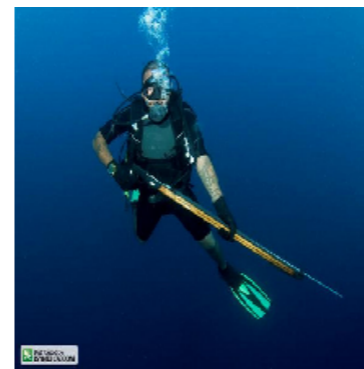
# Robustness Analysis of Neural Networks

**Goal:** identify adversarial inputs or check how amenable the network is for adversarial inputs

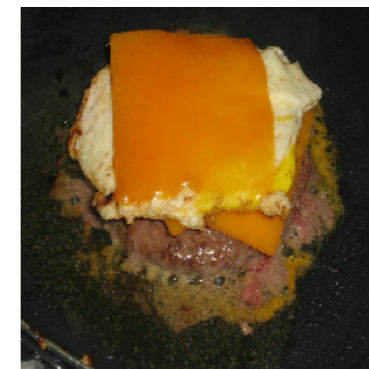
[Pei2017]

adversarial input

- **hardly** perceptible perturbations
- **large** impact on network's output



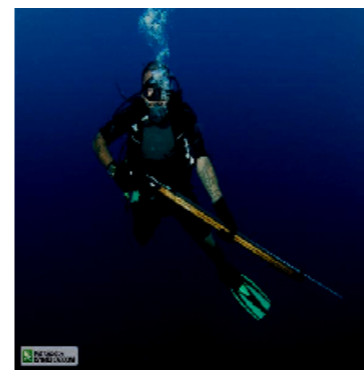
all:diver



all:cheeseburger



all:flamingo



IMG\_C1:ski



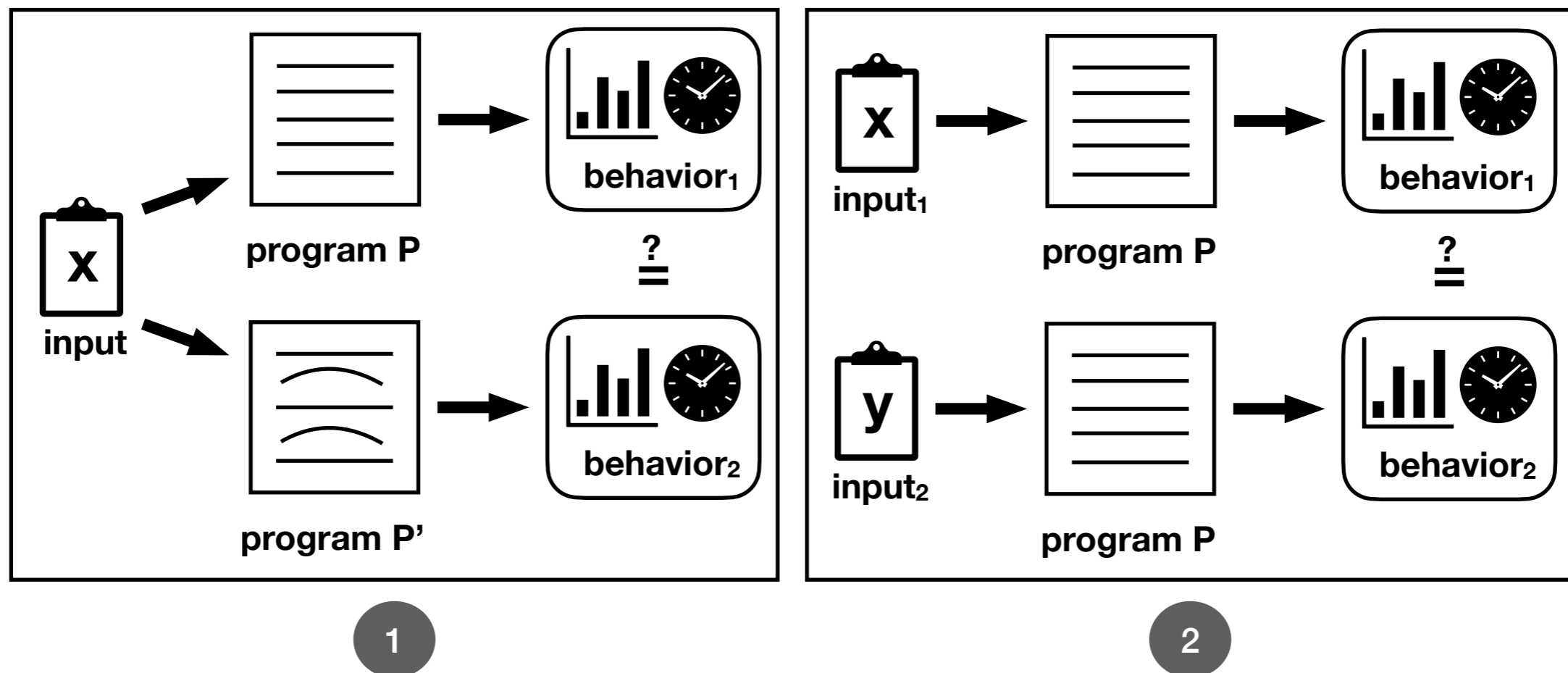
IMG\_C2:icecream



IMG\_C3:goldfish

# (My) Research Problem

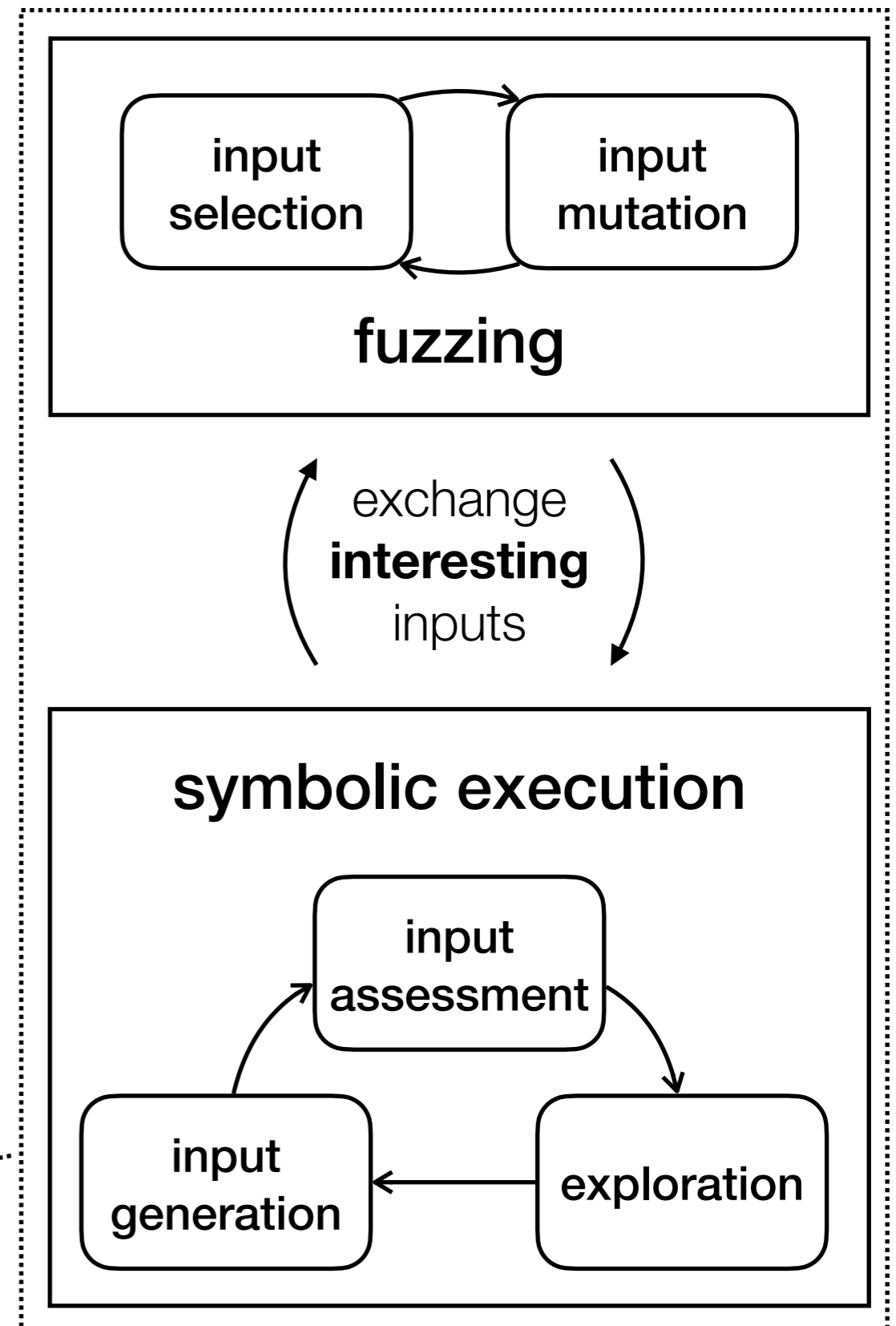
↪ identify behavioral differences



# Core Contributions

- (1) the concept of differential fuzzing
- (2) the concept of differential dynamic symbolic execution
- (3) the concept of hybrid analysis in differential program analysis
- (4) the concept of a hybrid setup for applying fuzzing and symbolic execution in parallel

**HyDiff**

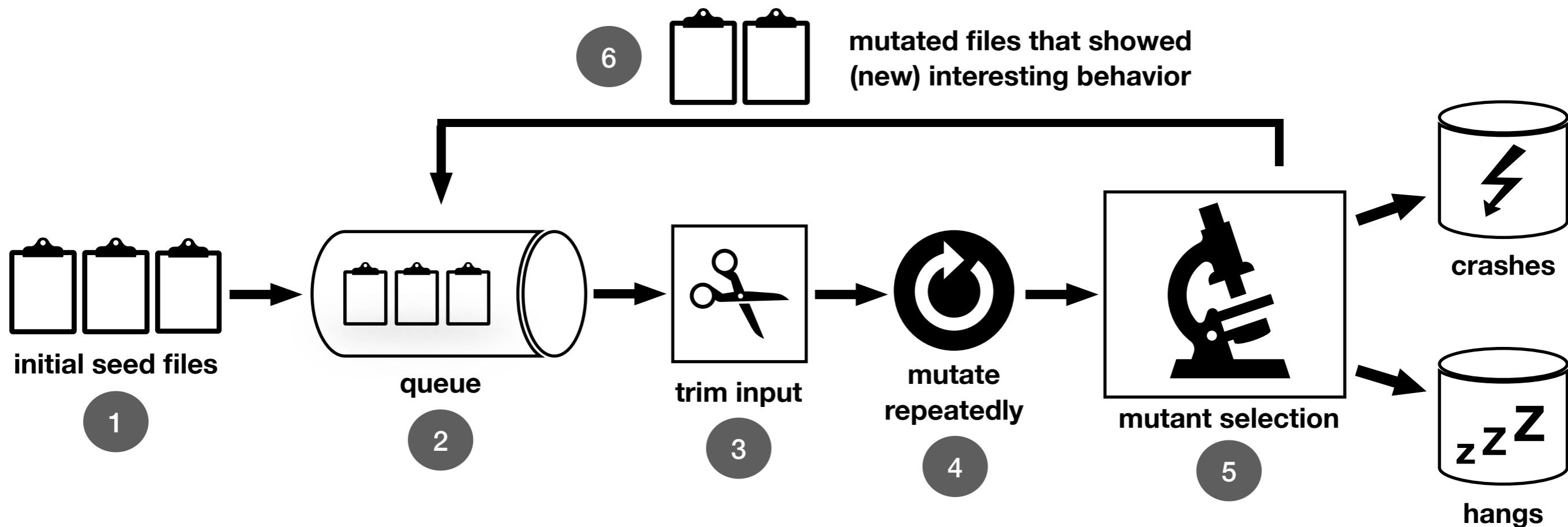


# Fuzzing

- term **fuzzing** was coined by Miller et al. in 1990, when they used a random testing tool to investigate the reliability of UNIX tools [Miller1990]
- classification based on degree of program analysis
  - blackbox / greybox / whitebox fuzzing
- classification based on generation technique
  - search-based fuzzing
  - generative fuzzing
- **state-of-the-art** in vulnerability detection: coverage-based, mutational fuzzing



# Coverage-Based Mutational Fuzzing



# Symbolic Execution

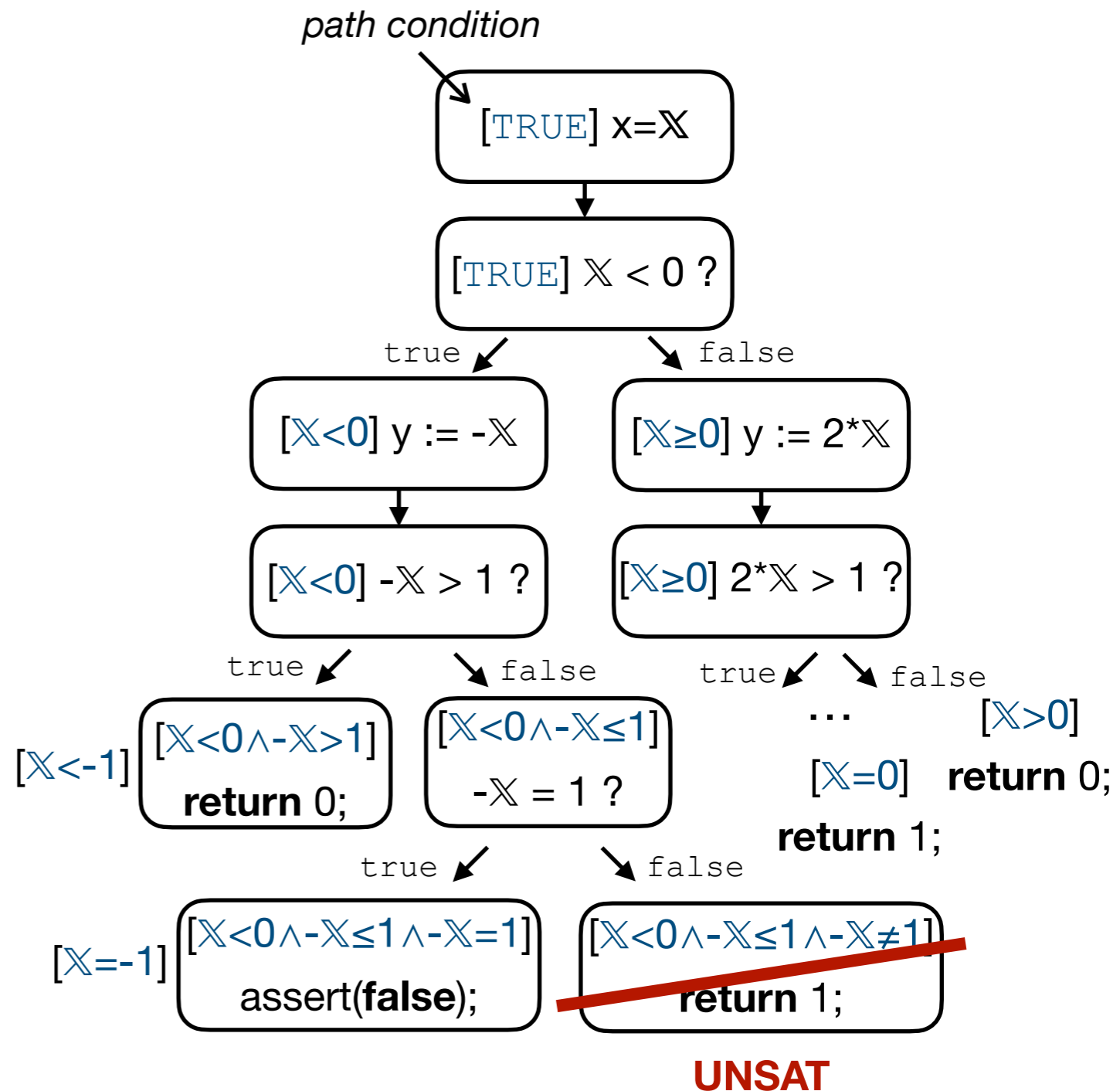
- introduced by King, Clarke, and Boyer et al. [King1976]  
[Clarke1976]  
[Boyer1975]
- analysis of programs with **unspecified inputs**, i.e.  
execute a program with **symbolic** inputs
- for each path, build a **path condition**

# Example: Symbolic Execution

```

1  int foo (int x) {
2      int y;
3      if (x < 0) {
4          y = -x;
5      } else {
6          y = 2 * x;
7      }
8      if (y > 1) {
9          return 0;
10     } else {
11         if (y == 1)
12             assert(false);
13     }
14     return 1;
15 }

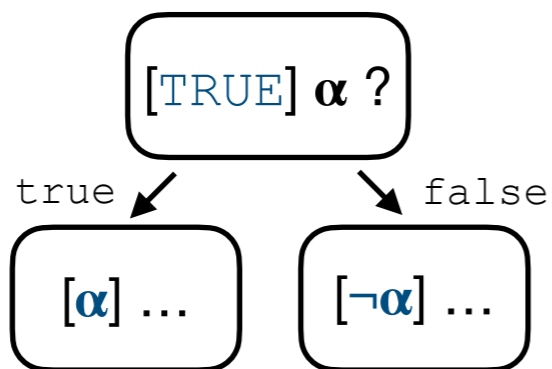
```



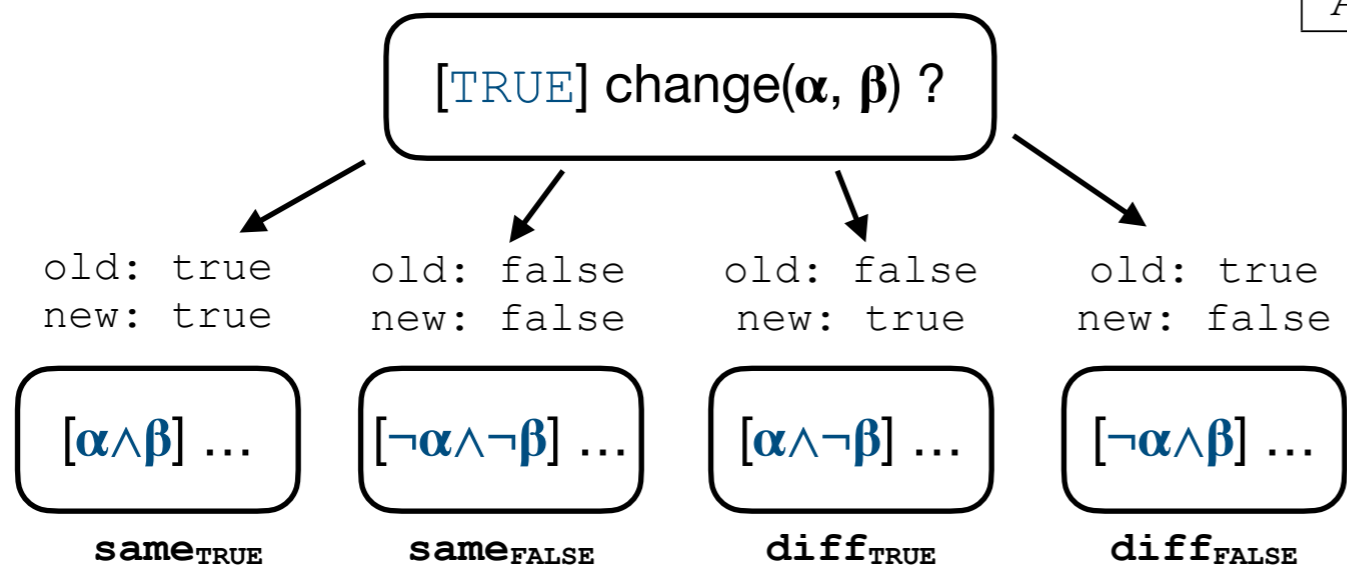
# Shadow Symbolic Execution

[Palikareva2016]

## Two-way Forking



## Four-way Forking



Change Type	Example
Update assignment	<code>x = x + <b>change</b>(E1, E2);</code>
Update condition	<code>if(<b>change</b>(E1, E2)) ...</code>
Add extra assignment	<code>x = <b>change</b>(x, E);</code>
Remove assignment	<code>x = <b>change</b>(E, x);</code>
Add conditional	<code>if(<b>change</b>(false, C)) ...</code>
Remove conditional	<code>if(<b>change</b>(C, false)) ...</code>
Remove code	<code>if(<b>change</b>(true, false)) ...</code>
Add code	<code>if(<b>change</b>(false, true)) ...</code>

# Why combine Fuzzing and Symbolic Execution?

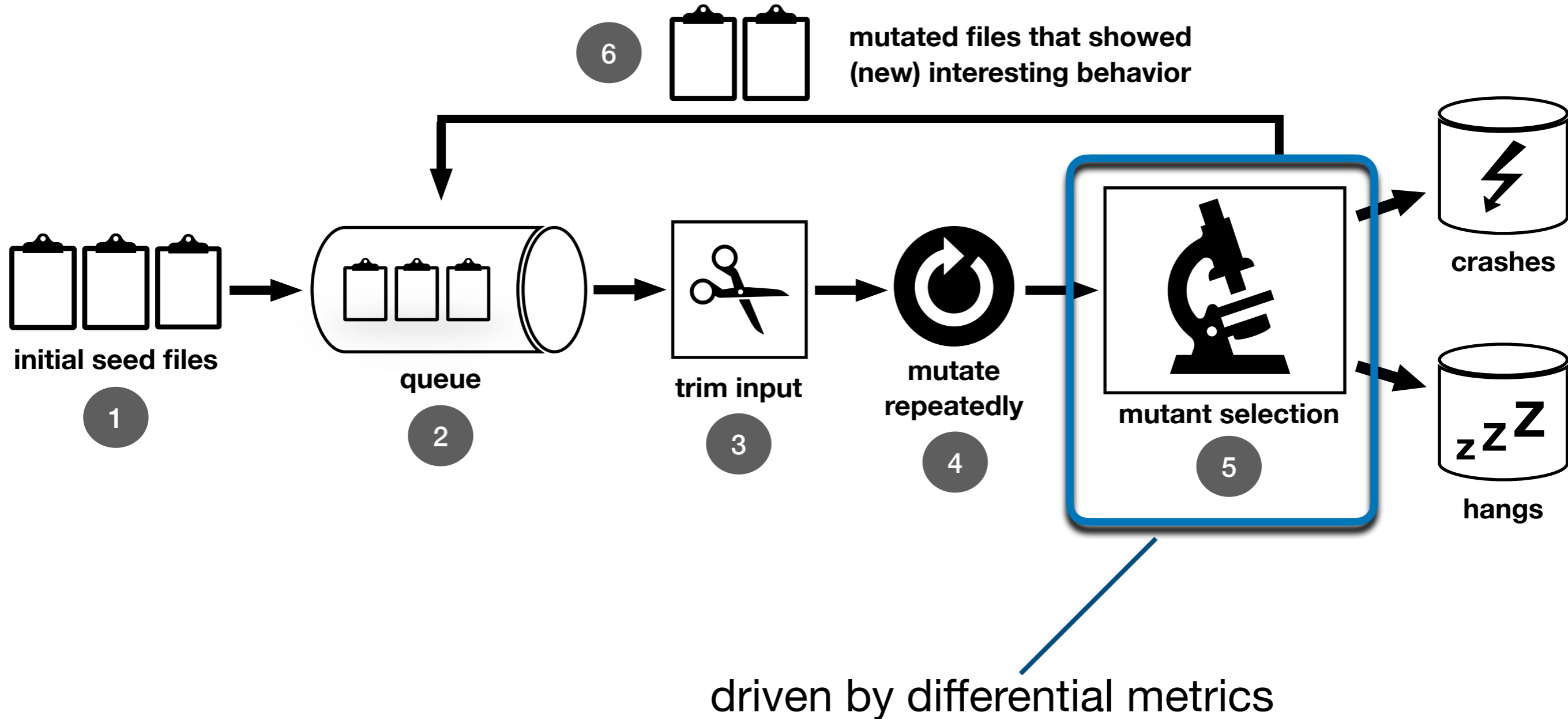
**good** in finding  
**shallow** bugs, but **bad**  
in finding **deep** program  
paths

input **reasoning ability**,  
but **path explosion** and  
**constraint solving**

# Related Work

- **regression analysis**  
[Person2008, Person2011, Yang2012, Orso2008, Taneja2008]
  - **side-channel analysis**  
[Antonopoulos2017, Chen2017, Pasareanu2016, Brennan2018]
  - **worst-case complexity analysis**  
[Petsios2017, Lemieux2018, Burnim2009, Luckow2017]
  - **robustness analysis of neural networks**  
[Ma2018, Pei2017, Sun2018, Goodfellow2014, Tian2018]
- ➔ **not directed to differential behavior**
  - ➔ **typical fuzzing problems**
  - ➔ **exhaustive exploration necessary**
  - ➔ **abstractions, bounded analysis, depend on models**

# Differential Fuzzing



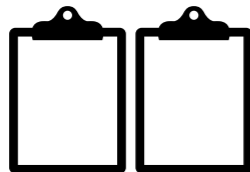
# Differential Metrics

- output difference (odiff)
- decision difference (ddiff)
- cost difference (cdiff)
- patch distance (only for regression testing)



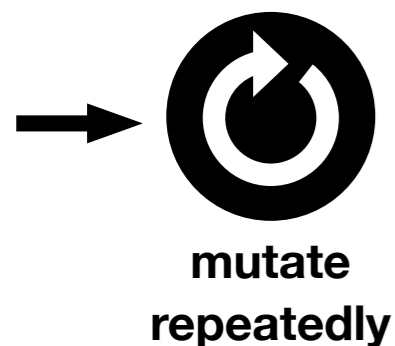
# Differential Fuzzing

6

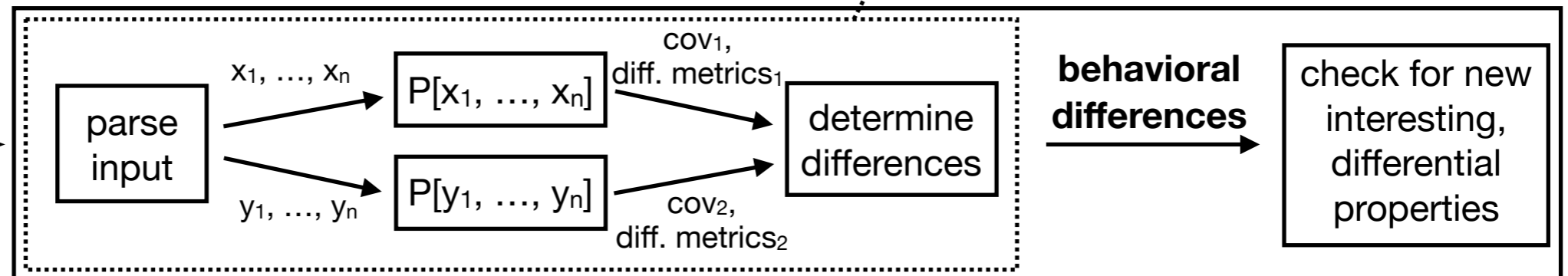


mutated files that showed  
(new) interesting behavior

fuzzing driver



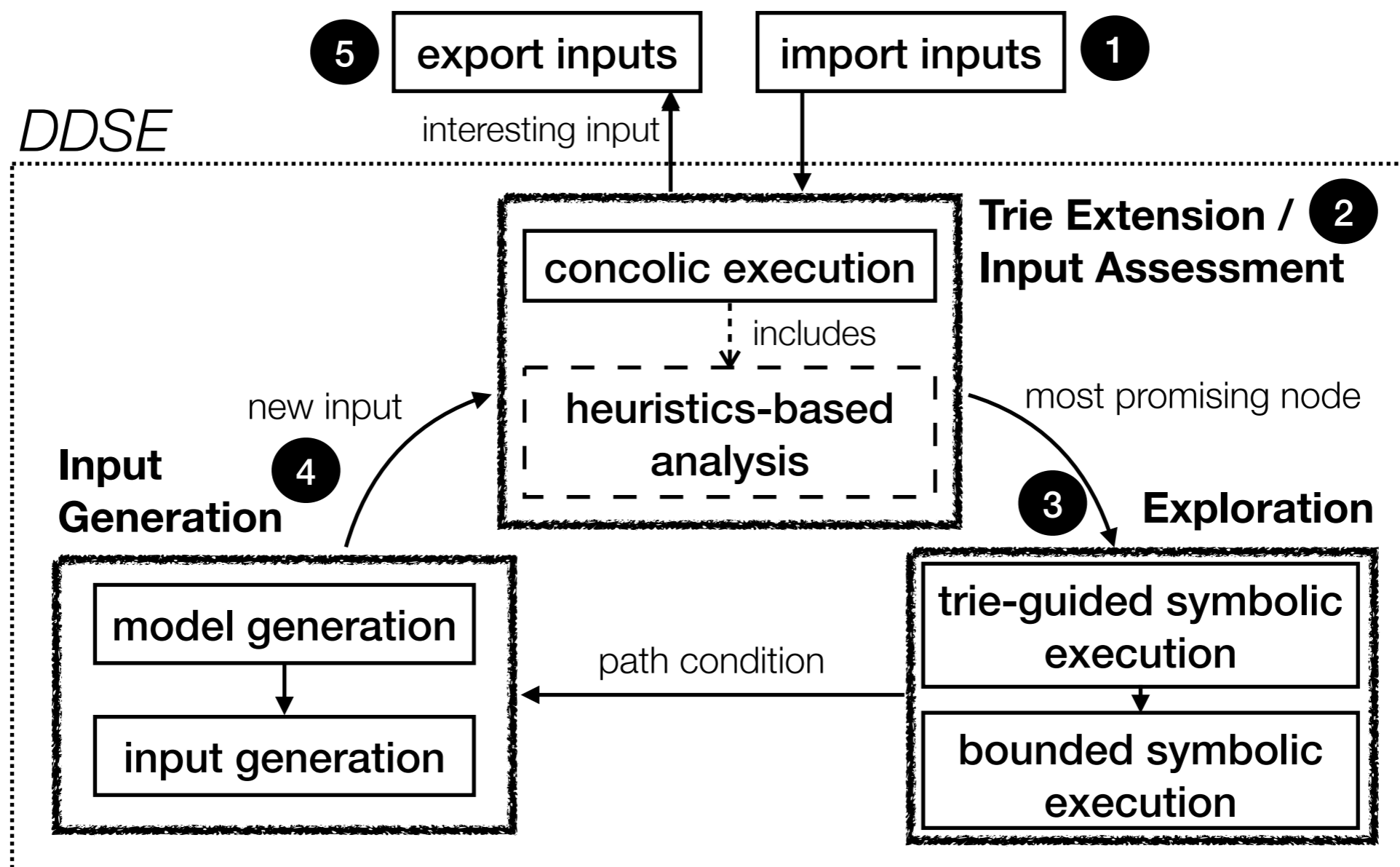
4



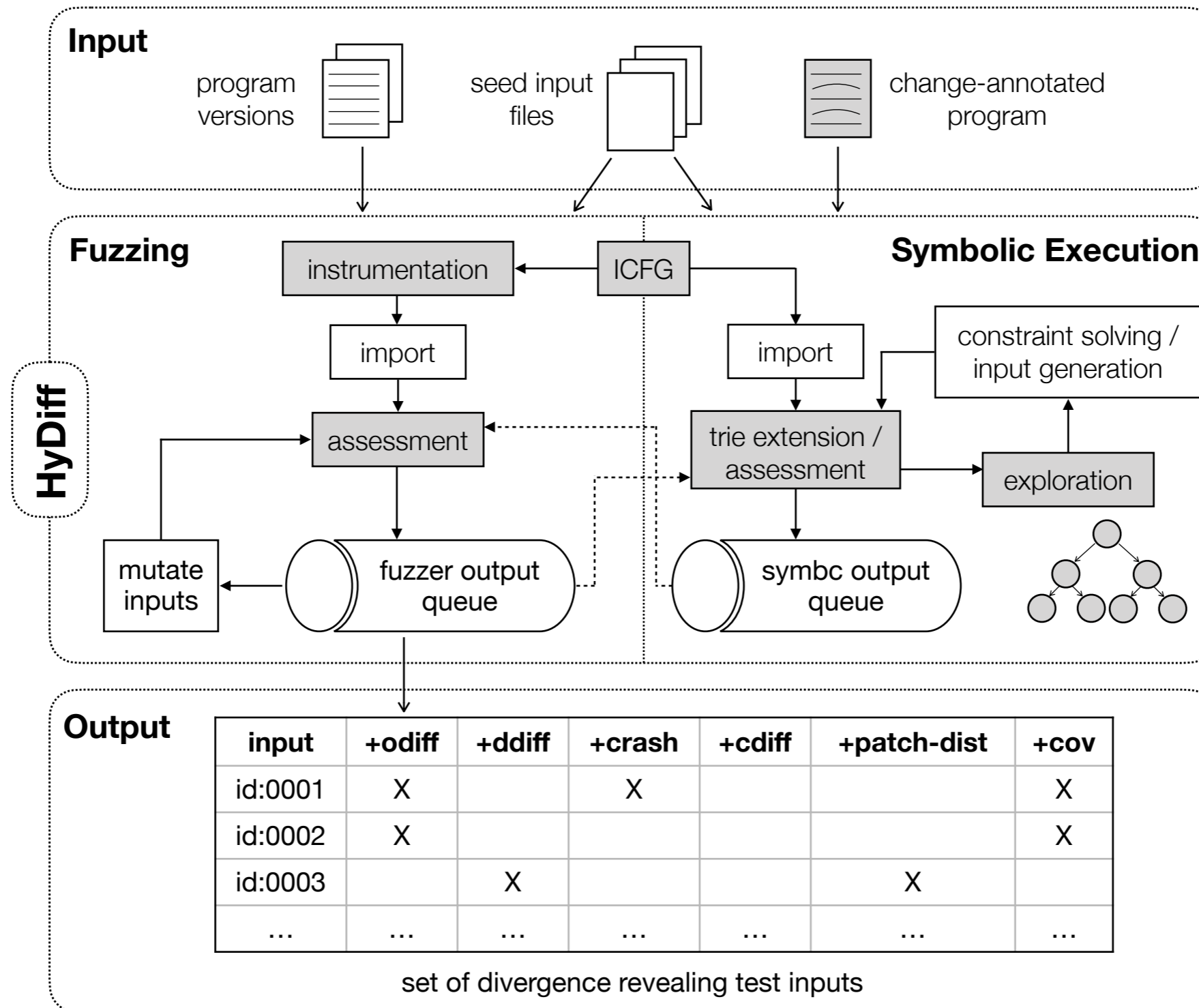
5

mutant selection by input evaluation for the  
instrumented program P

# Differential Dynamic SymExe



# HyDiff's overview



# Research Questions

**RQ1:** How good is *differential fuzzing* and what are the limitations?

**RQ2:** How good is *differential dynamic symbolic execution* and what are the limitations?

**RQ3:** Can the *hybrid* approach outperform the single techniques?

**RQ4:** Can the hybrid approach *not only combine* the results of fuzzing and symbolic execution, but also *amplify* the search itself and generate even better results than each approach on its own?

**RQ5:** Can the proposed hybrid differential software testing approach *reveal behavioral differences* in software?

# Evaluation Strategy

Quantitative analysis based on benchmarks in the specific application areas in differential analysis:

**A1**

Regression Analysis

**A2**

Worst-Case Complexity  
Analysis

**A3**

Side-Channel Analysis

**A4**

Robustness Analysis of  
Neural Networks

# Evaluation Metrics

**A1**

Regression Analysis

- average time to first output difference ( $t_{+odiff}$ )

**A4**

Robustness Analysis of Neural Networks

- $t_{min}$
- average output differences ( $\#odiff$ )
- average decision differences ( $\#ddiff$ )

**A2**

Worst-Case Complexity Analysis

- average maximum cost
- $cost_{max}$

**A3**

Side-Channel Analysis

- time to first cost improvement

# Evaluation Infrastructure

## What to compare?

Differential Fuzzing (DF)

**Parallel** Differential Fuzzing (PDF)

Differential Dynamic Symbolic Execution (DDSE)

DDSE with **double** time budget (DDSEx2T)

Hybrid Differential Software Testing (HyDiff)

# Regression Analysis

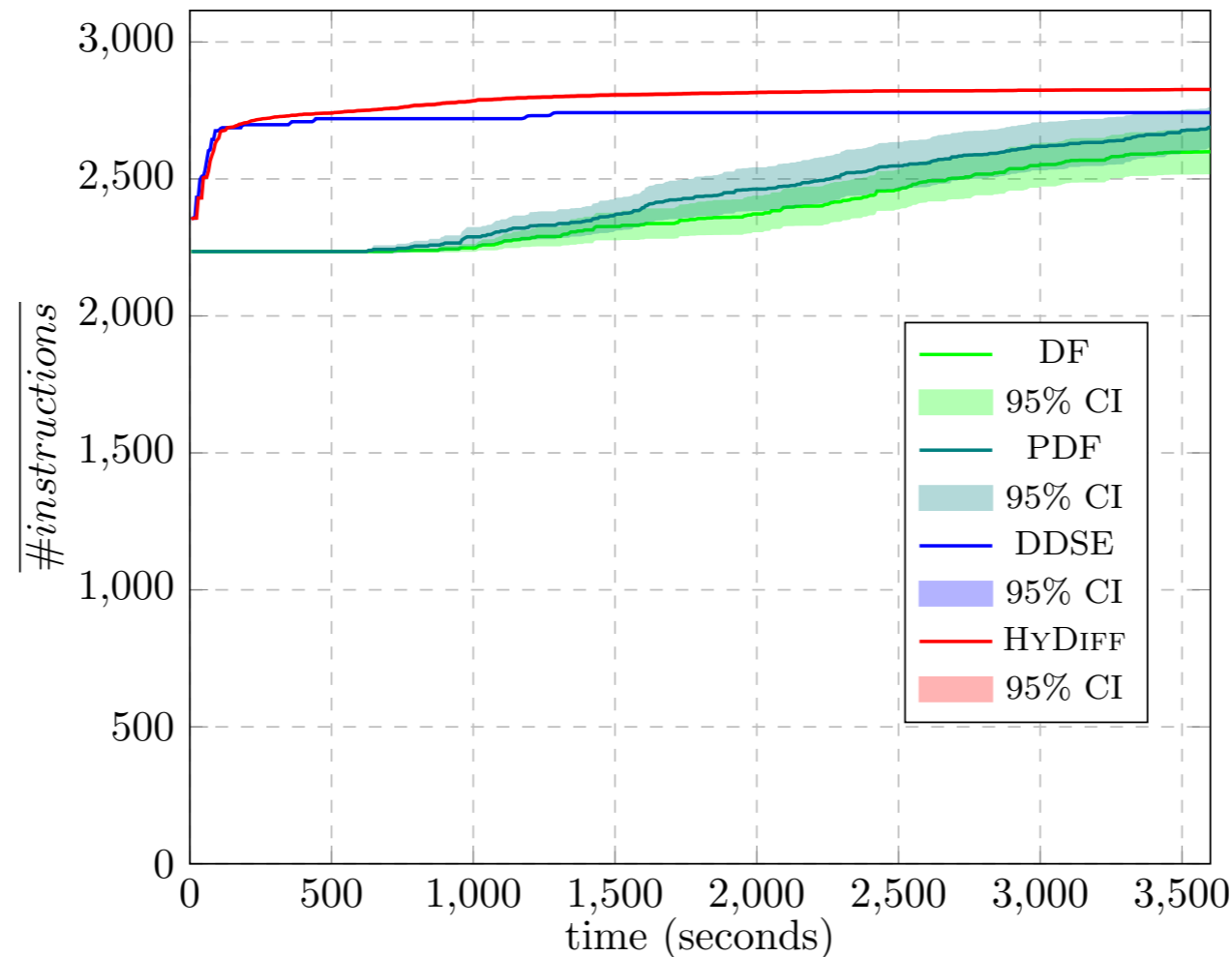
Subject (# changes)	Differential Fuzzing (DF)				Parallel Differential Fuzzing (PDF)				Differential Dynamic Sym. Exec. (DDSE)				DDSE double time budget (DDSEx2T)				HyDiff			
	$\bar{t}$	$t_{\min}$	#odiff	#ddiff	$\bar{t}$	$t_{\min}$	#odiff	#ddiff	$\bar{t}$	$t_{\min}$	#odiff	#ddiff	$\bar{t}$	$t_{\min}$	#odiff	#ddiff	$\bar{t}$	$t_{\min}$	#odiff	#ddiff
TCAS-1 (1)	-	-	0.00 ( $\pm 0.00$ )	0.00 ( $\pm 0.00$ )	-	-	0.00 ( $\pm 0.00$ )	0.00 ( $\pm 0.00$ )	<b>20.10</b> ( $\pm 0.14$ )	19	1.00 ( $\pm 0.00$ )	3.00 ( $\pm 0.00$ )	<b>20.10</b> ( $\pm 0.14$ )	19	1.00 ( $\pm 0.00$ )	3.00 ( $\pm 0.00$ )	49.87 ( $\pm 5.48$ )	29	1.00 ( $\pm 0.00$ )	<b>4.67</b> ( $\pm 0.40$ )
TCAS-2 (1)	441.83 ( $\pm 57.70$ )	120	0.70 ( $\pm 0.23$ )	2.13 ( $\pm 0.73$ )	335.93 ( $\pm 58.24$ )	16	1.57 ( $\pm 0.33$ )	5.40 ( $\pm 1.29$ )	<b>170.07</b> ( $\pm 0.32$ )	168	1.00 ( $\pm 0.00$ )	9.00 ( $\pm 0.00$ )	<b>170.07</b> ( $\pm 0.32$ )	168	1.00 ( $\pm 0.00$ )	9.00 ( $\pm 0.00$ )	186.87 ( $\pm 12.30$ )	92	1.23 ( $\pm 0.18$ )	<b>13.83</b> ( $\pm 0.37$ )
TCAS-3 (1)	588.43 ( $\pm 15.18$ )	392	0.10 ( $\pm 0.11$ )	38.63 ( $\pm 1.96$ )	531.87 ( $\pm 30.90$ )	295	0.67 ( $\pm 0.27$ )	55.53 ( $\pm 2.18$ )	<b>230.37</b> ( $\pm 0.52$ )	228	2.00 ( $\pm 0.00$ )	19.00 ( $\pm 0.00$ )	<b>230.37</b> ( $\pm 0.52$ )	228	2.00 ( $\pm 0.00$ )	19.00 ( $\pm 0.00$ )	263.20 ( $\pm 3.61$ )	236	2.00 ( $\pm 0.00$ )	57.43 ( $\pm 1.54$ )
TCAS-4 (1)	28.47 ( $\pm 10.42$ )	2	1.00 ( $\pm 0.00$ )	18.27 ( $\pm 1.06$ )	<b>9.27</b> ( $\pm 3.34$ )	1	1.00 ( $\pm 0.00$ )	24.10 ( $\pm 1.24$ )	-	-	0.00 ( $\pm 0.00$ )	3.00 ( $\pm 0.00$ )	-	-	0.00 ( $\pm 0.00$ )	3.00 ( $\pm 0.00$ )	43.70 ( $\pm 14.01$ )	3	1.00 ( $\pm 0.00$ )	22.53 ( $\pm 1.01$ )
TCAS-5 (1)	184.93 ( $\pm 46.66$ )	24	2.00 ( $\pm 0.00$ )	31.97 ( $\pm 1.06$ )	79.77 ( $\pm 21.40$ )	3	2.00 ( $\pm 0.00$ )	40.00 ( $\pm 1.73$ )	173.40 ( $\pm 0.34$ )	171	2.00 ( $\pm 0.00$ )	23.00 ( $\pm 0.00$ )	173.40 ( $\pm 0.34$ )	171	2.00 ( $\pm 0.00$ )	23.00 ( $\pm 0.00$ )	94.60 ( $\pm 30.72$ )	1	2.00 ( $\pm 0.00$ )	<b>49.83</b> ( $\pm 1.27$ )
TCAS-6 (1)	233.63 ( $\pm 54.48$ )	4	0.97 ( $\pm 0.06$ )	4.13 ( $\pm 0.83$ )	114.63 ( $\pm 37.12$ )	15	1.00 ( $\pm 0.00$ )	9.50 ( $\pm 0.98$ )	<b>4.73</b> ( $\pm 0.16$ )	4	1.00 ( $\pm 0.00$ )	6.00 ( $\pm 0.00$ )	<b>4.73</b> ( $\pm 0.16$ )	4	1.00 ( $\pm 0.00$ )	6.00 ( $\pm 0.00$ )	7.57 ( $\pm 0.26$ )	6	1.00 ( $\pm 0.00$ )	10.37 ( $\pm 0.70$ )
TCAS-7 (1)	-	-	0.00 ( $\pm 0.00$ )	0.00 ( $\pm 0.00$ )	581.60 ( $\pm 28.73$ )	164	0.07 ( $\pm 0.09$ )	0.27 ( $\pm 0.36$ )	73.50 ( $\pm 0.20$ )	72	2.00 ( $\pm 0.00$ )	6.00 ( $\pm 0.00$ )	73.50 ( $\pm 0.20$ )	72	2.00 ( $\pm 0.00$ )	6.00 ( $\pm 0.00$ )	<b>71.70</b> ( $\pm 1.71$ )	62	2.00 ( $\pm 0.00$ )	<b>8.93</b> ( $\pm 0.39$ )
TCAS-8 (1)	-	-	0.00 ( $\pm 0.00$ )	0.00 ( $\pm 0.00$ )	-	-	0.00 ( $\pm 0.00$ )	0.00 ( $\pm 0.00$ )	78.73 ( $\pm 1.24$ )	75	2.00 ( $\pm 0.00$ )	6.00 ( $\pm 0.00$ )	78.73 ( $\pm 1.24$ )	75	2.00 ( $\pm 0.00$ )	6.00 ( $\pm 0.00$ )	<b>65.33</b> ( $\pm 0.75$ )	61	2.00 ( $\pm 0.00$ )	<b>8.77</b> ( $\pm 0.49$ )
TCAS-9 (1)	221.73 ( $\pm 48.83$ )	10	1.00 ( $\pm 0.00$ )	6.13 ( $\pm 0.85$ )	<b>109.73</b> ( $\pm 28.35$ )	4	1.00 ( $\pm 0.00$ )	9.37 ( $\pm 0.44$ )	148.57 ( $\pm 1.76$ )	143	1.00 ( $\pm 0.00$ )	15.00 ( $\pm 0.00$ )	148.57 ( $\pm 1.76$ )	143	1.00 ( $\pm 0.00$ )	15.00 ( $\pm 0.00$ )	185.53 ( $\pm 18.42$ )	39	1.00 ( $\pm 0.00$ )	<b>22.37</b> ( $\pm 0.89$ )
TCAS-10 (2)	173.47 ( $\pm 46.27$ )	1	1.93 ( $\pm 0.09$ )	12.27 ( $\pm 1.69$ )	100.53 ( $\pm 25.20$ )	3	2.00 ( $\pm 0.00$ )	18.07 ( $\pm 1.07$ )	<b>4.87</b> ( $\pm 0.52$ )	4	2.00 ( $\pm 0.00$ )	12.00 ( $\pm 0.00$ )	<b>4.87</b> ( $\pm 0.52$ )	4	2.00 ( $\pm 0.00$ )	12.00 ( $\pm 0.00$ )	7.63 ( $\pm 0.22$ )	7	2.00 ( $\pm 0.00$ )	<b>21.30</b> ( $\pm 0.82$ )
Math-10 (1)	221.13 ( $\pm 56.26$ )	10	64.50 ( $\pm 15.98$ )	15.50 ( $\pm 2.35$ )	109.53 ( $\pm 18.08$ )	13	<b>172.37</b> ( $\pm 26.21$ )	24.03 ( $\pm 1.33$ )	<b>2.97</b> ( $\pm 0.17$ )	2	7.00 ( $\pm 0.00$ )	10.00 ( $\pm 0.00$ )	<b>2.97</b> ( $\pm 0.17$ )	2	7.00 ( $\pm 0.00$ )	10.00 ( $\pm 0.00$ )	3.87 ( $\pm 0.20$ )	3	44.33 ( $\pm 5.47$ )	<b>32.00</b> ( $\pm 1.39$ )
Math-46 (1)	377.87 ( $\pm 63.43$ )	77	0.80 ( $\pm 0.14$ )	36.33 ( $\pm 1.07$ )	270.07 ( $\pm 50.22$ )	8	1.00 ( $\pm 0.00$ )	<b>43.03</b> ( $\pm 0.78$ )	<b>118.93</b> ( $\pm 0.90$ )	116	1.00 ( $\pm 0.00$ )	5.60 ( $\pm 0.18$ )	<b>118.93</b> ( $\pm 0.90$ )	116	1.00 ( $\pm 0.00$ )	8.00 ( $\pm 0.00$ )	122.00 ( $\pm 8.34$ )	49	1.00 ( $\pm 0.00$ )	38.17 ( $\pm 0.82$ )
Math-60 (7)	6.93 ( $\pm 0.63$ )	4	219.17 ( $\pm 5.26$ )	92.90 ( $\pm 1.64$ )	5.90 ( $\pm 0.47$ )	4	<b>483.03</b> ( $\pm 9.52$ )	<b>138.10</b> ( $\pm 3.56$ )	<b>2.27</b> ( $\pm 0.16$ )	2	2.00 ( $\pm 0.00$ )	3.00 ( $\pm 0.00$ )	<b>2.27</b> ( $\pm 0.16$ )	2	2.00 ( $\pm 0.00$ )	3.00 ( $\pm 0.00$ )	4.77 ( $\pm 0.15$ )	4	234.23 ( $\pm 5.63$ )	94.20 ( $\pm 2.67$ )
Time-1 (14)	5.17 ( $\pm 1.20$ )	2	123.30 ( $\pm 5.86$ )	170.63 ( $\pm 3.43$ )	3.30 ( $\pm 0.60$ )	2	<b>221.00</b> ( $\pm 7.84$ )	<b>249.10</b> ( $\pm 4.29$ )	5.23 ( $\pm 0.18$ )	4	33.00 ( $\pm 0.00$ )	32.00 ( $\pm 0.00$ )	5.23 ( $\pm 0.18$ )	4	33.00 ( $\pm 0.00$ )	32.00 ( $\pm 0.00$ )	3.80 ( $\pm 0.69$ )	1	189.73 ( $\pm 11.94$ )	225.33 ( $\pm 5.62$ )
CLI1-2 (13)	-	-	0.00 ( $\pm 0.00$ )	159.53 ( $\pm 4.05$ )	-	-	0.00 ( $\pm 0.00$ )	<b>202.17</b> ( $\pm 3.48$ )	-	-	0.00 ( $\pm 0.00$ )	4.00 ( $\pm 0.00$ )	-	-	0.00 ( $\pm 0.00$ )	4.00 ( $\pm 0.00$ )	-	-	0.00 ( $\pm 0.00$ )	169.40 ( $\pm 4.07$ )
CLI2-3 (13)	10.83 ( $\pm 3.33$ )	2	82.30 ( $\pm 3.98$ )	176.83 ( $\pm 3.62$ )	<b>4.83</b> ( $\pm 1.29$ )	1	<b>161.60</b> ( $\pm 6.62$ )	242.53 ( $\pm 6.92$ )	-	-	0.00 ( $\pm 0.00$ )	37.00 ( $\pm 0.00$ )	-	-	0.00 ( $\pm 0.00$ )	37.00 ( $\pm 0.00$ )	13.27 ( $\pm 3.62$ )	2	84.63 ( $\pm 4.24$ )	242.70 ( $\pm 3.80$ )
CLI3-4 (8)	7.43 ( $\pm 1.60$ )	1	96.73 ( $\pm 4.54$ )	279.13 ( $\pm 4.51$ )	7.20 ( $\pm 1.85$ )	2	97.87 ( $\pm 4.02$ )	467.27 ( $\pm 5.05$ )	<b>4.07</b> ( $\pm 0.36$ )	3	1.00 ( $\pm 0.00$ )	12.00 ( $\pm 0.00$ )	<b>4.07</b> ( $\pm 0.36$ )	3	1.00 ( $\pm 0.00$ )	12.00 ( $\pm 0.00$ )	8.93 ( $\pm 2.13$ )	2	<b>113.33</b> ( $\pm 4.80$ )	471.50 ( $\pm 8.93$ )
CLI4-5 (13)	589.57 ( $\pm 16.05$ )	358	0.07 ( $\pm 0.09$ )	219.30 ( $\pm 3.74$ )	-	-	0.00 ( $\pm 0.00$ )	<b>274.43</b> ( $\pm 4.22$ )	-	-	0.00 ( $\pm 0.00$ )	4.00 ( $\pm 0.00$ )	-	-	0.00 ( $\pm 0.00$ )	4.00 ( $\pm 0.00$ )	<b>551.97</b> ( $\pm 45.65$ )	125	<b>0.13</b> ( $\pm 0.12$ )	235.17 ( $\pm 5.73$ )
CLI5-6 (21)	<b>4.13</b> ( $\pm 1.04$ )	1	143.87 ( $\pm 4.99$ )	182.00 ( $\pm 5.54$ )	<b>3.43</b> ( $\pm 0.72$ )	1	<b>277.17</b> ( $\pm 6.81$ )	<b>272.17</b> ( $\pm 7.32$ )	-	-	0.00 ( $\pm 0.00$ )	5.00 ( $\pm 0.00$ )	-	-	0.00 ( $\pm 0.00$ )	5.00 ( $\pm 0.00$ )	6.17 ( $\pm 1.31$ )	2	177.80 ( $\pm 4.39$ )	214.47 ( $\pm 6.38$ )

HyDiff classifies all subjects correctly.

Components do benefit from each other.



# Worst-Case Complexity Analysis



*Hexcolor*

DDSE quickly makes progress, DF continuously improves the score.

HyDiff successfully combines strengths of DDSE and DF.

HyDiff outperforms the components in isolation.

# Side-Channel Analysis

- in regression testing: **changes** in the **program**
- in side-channel analysis: **changes** in the **input**

`secret = change(secret1, secret2)`

# Side-Channel Analysis

Benchmark	Version	Differential Fuzzing (DF)			Themis		
		$\bar{\delta}$	$\delta_{\max}$	$\bar{t} : \delta > 0$	$\epsilon = 64$	$\epsilon = 0$	Time (s)
Spring-Security	Safe	1.00 ( $\pm 0.00$ )	1	4.77 ( $\pm 1.07$ )	✓	✓	1.70
Spring-Security	Unsafe	149.00 ( $\pm 0.00$ )	149	4.17 ( $\pm 0.90$ )	✓	✓	1.09
JDK7-MsgDigest	Safe	1.00 ( $\pm 0.00$ )	1	10.77 ( $\pm 2.12$ )	✓	✓	1.27
JDK6-MsgDigest	Unsafe	140.03 ( $\pm 20.39$ )	263	3.20 ( $\pm 0.81$ )	✓	✓	1.33
Picketbox	Safe	1.00 ( $\pm 0.00$ )	1	16.90 ( $\pm 3.89$ )	✓	✗	1.79
Picketbox	Unsafe	363.70 ( $\pm 562.18$ )	8,822	5.13 ( $\pm 1.83$ )	✓	✓	1.55
Tomcat	Safe	25.07 ( $\pm 0.36$ )	26	19.90 ( $\pm 9.29$ )	✓	✗	9.93
<i>Tomcat</i>	<i>Unsafe</i>	<i>49.00 (<math>\pm 0.36</math>)</i>	<i>50</i>	<i>23.53 (<math>\pm 9.73</math>)</i>	<i>✓</i>	<i>✓</i>	<i>8.64</i>
<i>Jetty</i>	<i>Safe</i>	<i>11.77 (<math>\pm 0.60</math>)</i>	<i>15</i>	<i>3.77 (<math>\pm 0.72</math>)</i>	<i>✓</i>	<i>✓</i>	<i>2.50</i>
Jetty	Unsafe	70.87 ( $\pm 6.12$ )	105	6.83 ( $\pm 1.62$ )	✓	✓	2.07
orientdb	Safe	1.00 ( $\pm 0.00$ )	1	16.60 ( $\pm 5.14$ )	✓	✗	37.99
orientdb	Unsafe	458.93 ( $\pm 685.64$ )	10,776	4.77 ( $\pm 1.06$ )	✓	✓	38.09
pac4j	Safe	10.00 ( $\pm 0.00$ )	10	1.10 ( $\pm 0.11$ )	✓	✗	3.97
<i>pac4j</i>	<i>Unsafe</i>	<i>11.00 (<math>\pm 0.00</math>)</i>	<i>11</i>	<i>1.13 (<math>\pm 0.12</math>)</i>	<i>✓</i>	<i>✓</i>	<i>1.85</i>
<i>pac4j</i>	<i>Unsafe*</i>	<i>39.00 (<math>\pm 0.00</math>)</i>	<i>39</i>	<i>1.10 (<math>\pm 0.11</math>)</i>	<i>-</i>	<i>-</i>	<i>-</i>

DF can find the same vulnerabilities as static analysis

well-balanced combination: fast and high delta  
(important to assess the severity of vulnerability)

# Robustness Analysis of Neural Networks

**Purpose: stress test proposed technique**

- similar to SC analysis: **changes** in the **input**
- similar to regression analysis: **search for output differences**
- idea: allow up to x% changes in the pixels of the input image

$a[i][j] = \text{change}(a[i][j], \text{value});$

# NN Analysis

Subject (% change)	Differential Fuzzing (DF)				Parallel Differential Fuzzing (PDF)				Differential Dynamic Sym. Exec. (DDSE)				HyDiff			
	$\bar{t}$ +odiff	$t_{\min}$	$\overline{\#odiff}$	$\overline{\#ddiff}$	$\bar{t}$ +odiff	$t_{\min}$	$\overline{\#odiff}$	$\overline{\#ddiff}$	$\bar{t}$ +odiff	$t_{\min}$	$\overline{\#odiff}$	$\overline{\#ddiff}$	$\bar{t}$ +odiff	$t_{\min}$	$\overline{\#odiff}$	$\overline{\#ddiff}$
1	2,725.40 (+341.09)	1,074	0.57 ( $\pm 0.20$ )	7.73 ( $\pm 0.18$ )	2,928.60 ( $\pm 289.44$ )	1,202	1.00 ( $\pm 0.31$ )	<b>12.00</b> ( $\pm 0.48$ )	296.03 ( $\pm 1.49$ )	289	<b>1.00</b> ( $\pm 0.00$ )	1.00 ( $\pm 0.00$ )	297.10 ( $\pm 2.38$ )	267	<b>1.20</b> ( $\pm 0.14$ )	6.10 ( $\pm 0.11$ )
2	2,581.47 (+326.21)	1,032	0.93 ( $\pm 0.28$ )	7.93 ( $\pm 0.13$ )	2,509.20 ( $\pm 289.37$ )	1,117	1.23 ( $\pm 0.33$ )	<b>12.63</b> ( $\pm 0.48$ )	309.77 ( $\pm 7.04$ )	293	<b>1.00</b> ( $\pm 0.00$ )	1.00 ( $\pm 0.00$ )	<b>297.93</b> ( $\pm 1.29$ )	292	1.53 ( $\pm 0.20$ )	6.93 ( $\pm 0.13$ )
5	2,402.97 (+329.59)	1,189	1.23 ( $\pm 0.37$ )	6.47 ( $\pm 0.18$ )	2,501.43 ( $\pm 285.86$ )	1,429	1.70 ( $\pm 0.44$ )	<b>10.33</b> ( $\pm 0.43$ )	304.53 ( $\pm 1.06$ )	300	<b>1.00</b> ( $\pm 0.00$ )	1.00 ( $\pm 0.00$ )	<b>301.83</b> ( $\pm 1.16$ )	296	2.07 ( $\pm 0.29$ )	6.90 ( $\pm 0.17$ )
10	2,155.40 (+343.76)	996	1.57 ( $\pm 0.34$ )	8.10 ( $\pm 0.17$ )	2,127.70 ( $\pm 229.21$ )	1,418	2.20 ( $\pm 0.33$ )	<b>11.23</b> ( $\pm 0.40$ )	311.90 ( $\pm 0.74$ )	308	<b>1.00</b> ( $\pm 0.00$ )	1.00 ( $\pm 0.00$ )	311.07 ( $\pm 1.01$ )	306	2.37 ( $\pm 0.31$ )	7.00 ( $\pm 0.13$ )
20	1,695.83 (+228.18)	953	2.70 ( $\pm 0.37$ )	9.13 ( $\pm 0.12$ )	1,897.67 ( $\pm 219.97$ )	1,340	3.30 ( $\pm 0.49$ )	<b>11.57</b> ( $\pm 0.50$ )	346.87 ( $\pm 1.98$ )	339	<b>1.00</b> ( $\pm 0.00$ )	1.00 ( $\pm 0.00$ )	<b>341.83</b> ( $\pm 1.27$ )	336	3.13 ( $\pm 0.34$ )	7.20 ( $\pm 0.14$ )
50	1,830.83 (+259.79)	1,220	2.43 ( $\pm 0.42$ )	6.33 ( $\pm 0.21$ )	1,696.10 ( $\pm 86.20$ )	1,423	3.80 ( $\pm 0.35$ )	<b>12.00</b> ( $\pm 0.39$ )	455.03 ( $\pm 1.62$ )	449	<b>1.00</b> ( $\pm 0.00$ )	1.00 ( $\pm 0.00$ )	452.63 ( $\pm 2.06$ )	434	3.77 ( $\pm 0.34$ )	7.27 ( $\pm 0.16$ )
100	1,479.17 (+231.25)	960	2.47 ( $\pm 0.37$ )	9.37 ( $\pm 0.20$ )	1,790.87 ( $\pm 270.10$ )	1,109	3.03 ( $\pm 0.54$ )	<b>13.97</b> ( $\pm 0.68$ )	583.33 ( $\pm 2.83$ )	571	<b>1.00</b> ( $\pm 0.00$ )	1.00 ( $\pm 0.00$ )	<b>575.13</b> ( $\pm 2.65$ )	564	3.10 ( $\pm 0.35$ )	7.60 ( $\pm 0.18$ )

Shows the limitations of both components.

HyDiff can combine them so that both can benefit from each other

# RQ 1: Differential Fuzzing

- **Regression:** performs quite reasonable, but not all subject correctly classified (parallel DF did not help)
- **WCA:** improves cost continuously over time
- **SC:** outperforms Blazer and Themis
- **NN:** effective, but very slow (gets better with more x%)

Differential Fuzzing **continuously improves** its differential analysis over time

Parallel Differential Fuzzing **even better**, sometimes outperformed hybrid combination

# RQ 2: Differential Dynamic Symbolic Execution

- **Regression**: fast in finding output differences, but not all subject correctly classified
- **WCA**: often stays in plateaus without improvement, but good in finding some first slowdown
- **SC**: slow in the beginning, but eventually high delta
- **NN**: very fast for first output difference, but limited by heavy constraint solving

**DDSE develops in jumps** and only rarely in continuous improvement

**effective** technique due to **constraint solving**

**DDSE with twice the time budget does not improve the result**

# RQ 3+4: Hybrid combination

- **Regression:** HyDiff finds all output differences and often generates higher values in a shorter time period
- **WCA:** clearly outperforms components
- **SC:** no clear improvement, but well balanced combination
- **NN:** good combination, finds output differences and is fast

HyDiff **does not only combine results** of components but also **amplifies** them



# RQ 5: HyDiff for Differential Testing

- **Regression:** crashes not present, but inputs for behavioral differences
- **WCA:** AC vulnerabilities identified
- **SC:** all vulnerabilities identified
- **NN:** limits of HyDiff, however found adversarial inputs

HyDiff is **effective** for differential testing

# Publications

## **Shadow Symbolic Execution with Java PathFinder**

*Yannic Noller, Hoang Lam Nguyen, Minxing Tang, and Timo Kehrer*

Java Pathfinder Workshop 2017, SIGSOFT Software Engineering Notes 42 (January 2018)

## **Badger: Complexity Analysis with Fuzzing and Symbolic Execution**

*Yannic Noller, Rody Kersten, and Corina S. Păsăreanu*

ACM SIGSOFT International Symposium on Software Testing and Analysis (ISSTA) 2018

## **Differential Program Analysis with Fuzzing and Symbolic execution**

*Yannic Noller* (Doctoral Symposium Paper)

ACM/IEEE International Conference on Automated Software Engineering (ASE) 2018

## **DifFuzz: Differential Fuzzing for Side-Channel Analysis**

*Shirin Nilizadeh\*, Yannic Noller\*, and Corina S. Păsăreanu (\* joint first authors)*

ACM/IEEE International Conference on Software Engineering (ICSE) 2019

## **Complete Shadow Symbolic Execution with Java PathFinder**

*Yannic Noller, Hoang Lam Nguyen, Minxing Tang, Timo Kehrer and Lars Grunske*

Java Pathfinder Workshop 2019, SIGSOFT Software Engineering Notes 44 (December 2019)

## **HyDiff: Hybrid Differential Software Analysis**

*Yannic Noller, Corina S. Păsăreanu, Marcel Böhme, Youcheng Sun,*

*Hoang Lam Nguyen, and Lars Grunske*

ACM/IEEE International Conference on Software Engineering (ICSE) 2020

# Hybrid Differential Software Testing

Problem Contribution Background Solutions Validation Summary

## Differential Software Testing

1

2

yannic.noller@acm.org Hybrid Differential Software Testing 5

Problem Contribution Background Solutions Validation Summary

## Why combine Fuzzing and Symbolic Execution?

good in finding shallow bugs, but bad in finding deep program paths

input reasoning ability, but path explosion and constraint solving

yannic.noller@acm.org Hybrid Differential Software Testing 26

Problem Contribution Background Solutions Validation Summary

## Differential Fuzzing

4

5

6

yannic.noller@acm.org Hybrid Differential Software Testing 41

Problem Contribution Background Solutions Validation Summary

## Differential Dynamic SymExe

1

2

3

4

5

DDSE

yannic.noller@acm.org Hybrid Differential Software Testing 43

Problem Contribution Background Solutions Validation Summary

## HyDiff's overview

HyDiff

Input	+odiff	+ddiff	+crash	+cdiff	+patch-dist	+cov
id:0001	X		X			X
id:0002	X					X
id:0003		X			X	
...	...	...	...	...	...	...

set of divergence revealing test inputs

yannic.noller@acm.org Hybrid Differential Software Testing 50

Problem Contribution Background Solutions Validation Summary

## Regression Analysis TCAS-7

#diff/F

time (seconds)

Components do benefit from each other.

yannic.noller@acm.org Hybrid Differential Software Testing 36



# References

- [Antonopoulos2017]** Timos Antonopoulos, Paul Gazzillo, Michael Hicks, Eric Koskinen, Tachio Terauchi, and Shiyi Wei. 2017. Decomposition instead of self-composition for proving the absence of timing channels. In Proceedings of the 38th ACM SIGPLAN Conference on Programming Language Design and Implementation (PLDI 2017). ACM, New York, NY, USA, 362-375.
- [Barthe2004]** G. Barthe, P. R. D'Argenio and T. Rezk, "Secure Information Flow by Self-Composition," Computer Security Foundations Workshop, IEEE(CSFW), Pacific Grove, California, 2004, pp. 100.
- [Boyer1975]** Boyer, R. S., Elspas, B., & Levitt, K. N. (1975). SELECT - a Formal System for Testing and Debugging Programs by Symbolic Execution. SIGPLAN Not., 10(6), 234–245.
- [Brennan2018]** Brennan, T., Saha, S., Bultan, T., & Pasareanu, C. S. (2018). Symbolic Path Cost Analysis for Side-channel Detection. Proceedings of the 27th ACM SIGSOFT International Symposium on Software Testing and Analysis, 27–37.
- [Burnim2009]** J. Burnim, S. Juvekar, and K. Sen. 2009. WISE: Automated test generation for worst-case complexity. In 2009 IEEE 31st International Conference on Software Engineering. 463–473.
- [Cha2012]** Cha, S. K., Avgerinos, T., Rebert, A., & Brumley, D. (2012). Unleashing Mayhem on Binary Code. 2012 IEEE Symposium on Security and Privacy, 380–394.
- [Chen2017]** Jia Chen, Yu Feng, and Isil Dillig. 2017. Precise Detection of Side-Channel Vulnerabilities using Quantitative Cartesian Hoare Logic. In Proceedings of the 2017 ACM SIGSAC Conference on Computer and Communications Security (CCS '17). ACM, New York, NY, USA, 875-890.
- [Clarke1976]** L. A. Clarke, "A System to Generate Test Data and Symbolically Execute Programs," in IEEE Transactions on Software Engineering, vol. SE-2, no. 3, pp. 215-222, Sept. 1976.

# References (continued)

**[Goodfellow2014]** Goodfellow, I., Pouget-Abadie, J., Mirza, M., Xu, B., Warde-Farley, D., Ozair, S., Courville, A., & Bengio, Y. (2014). Generative Adversarial Nets. In Z. Ghahramani, M. Welling, C. Cortes, N. D. Lawrence, & K. Q. Weinberger (Eds.), *Advances in Neural Information Processing Systems 27* (pp. 2672–2680). Curran Associates, Inc.

**[IEEE2017]** ISO/IEC/IEEE International Standard - Systems and software engineering - Vocabulary. (2017). ISO/IEC/IEEE 24765:2017(E), 1–541.

**[King1976]** James C. King. 1976. Symbolic execution and program testing. *Commun. ACM* 19, 7 (July 1976), 385-394.

**[Lemieux2018]** Caroline Lemieux, Rohan Padhye, Koushik Sen, and Dawn Song. 2018. PerfFuzz: automatically generating pathological inputs. In *Proceedings of the 27th ACM SIGSOFT International Symposium on Software Testing and Analysis (ISSTA 2018)*. ACM, New York, NY, USA, 254-265.

**[Luckow2017]** Kasper Luckow, Rody Kersten, and Corina Pasareanu. 2017. Symbolic Complexity Analysis using Context-preserving Histories. In *Proceedings of the 10th IEEE International Conference on Software Testing, Verification and Validation (ICST 2017)*. 58–68

**[Ma2018]** Ma, L., Zhang, F., Sun, J., Xue, M., Li, B., Juefei-Xu, F., Xie, C., Li, L., Liu, Y., Zhao, J., & Wang, Y. (2018). DeepMutation: Mutation Testing of Deep Learning Systems. *2018 IEEE 29th International Symposium on Software Reliability Engineering (ISSRE)*, 100–111.

**[Miller1990]** Barton P. Miller, Louis Fredriksen, and Bryan So. 1990. An empirical study of the reliability of UNIX utilities. *Commun. ACM* 33, 12 (December 1990), 32-44.

# References (continued)

**[Ognawala2019]** Ognawala, S., Kilger, F., & Pretschner, A. (2019). Compositional Fuzzing Aided by Targeted Symbolic Execution. ArXiv Preprint ArXiv:1903.02981

**[Orso2008]** Orso, A., & Xie, T. (2008). BERT: BEhavioral Regression Testing. Proceedings of the 2008 International Workshop on Dynamic Analysis: Held in Conjunction with the ACM SIGSOFT International Symposium on Software Testing and Analysis (ISSTA 2008), 36–42.

**[Orso2014]** Orso, A., & Rothermel, G. (2014). Software Testing: A Research Travelogue (2000-2014). Proceedings of the on Future of Software Engineering, 117–132.

**[Pak2012]** Pak, B. S. (2012). Hybrid Fuzz Testing: Discovering Software Bugs via Fuzzing and Symbolic Execution. Carnegie Mellon University.

**[Palikareva2016]** Hristina Palikareva, Tomasz Kuchta, and Cristian Cadar. 2016. Shadow of a doubt: testing for divergences between software versions. In Proceedings of the 38th International Conference on Software Engineering (ICSE '16). ACM, New York, NY, USA, 1181-1192.

**[Pasareanu2016]** Pasareanu, C. S., Phan, Q.-S., & Malacaria, P. (2016). Multi-run Side-Channel Analysis Using Symbolic Execution and Max-SMT. 2016 IEEE 29th Computer Security Foundations Symposium (CSF), 387–400.

**[Petsios2017]** Petsios, T., Tang, A., Stolfo, S., Keromytis, A. D., & Jana, S. (2017). NEZHA: Efficient Domain-Independent Differential Testing. 2017 IEEE Symposium on Security and Privacy (SP), 615–632.

**[Pei2017]** Pei, K., Cao, Y., Yang, J., & Jana, S. (2017). DeepXplore: Automated Whitebox Testing of Deep Learning Systems. Proceedings of the 26th Symposium on Operating Systems Principles, 1–18.

# References (continued)

- [Person2008]** Person, S., Dwyer, M. B., Elbaum, S., & Pasareanu, C. S. (2008). Differential Symbolic Execution. Proceedings of the 16th ACM SIGSOFT International Symposium on Foundations of Software Engineering, 226–237.
- [Person2011]** Suzette Person, Guowei Yang, Neha Rungta, and Sarfraz Khurshid. 2011. Directed incremental symbolic execution. In Proceedings of the 32nd ACM SIGPLAN Conference on Programming Language Design and Implementation (PLDI '11). ACM, New York, NY, USA, 504-515.
- [Stephens2016]** Stephens N, Grosen J, Salls C, Dutcher A, Wang R, Corbetta J, Shoshitaishvili Y, Kruegel C, Vigna G. Driller: Augmenting Fuzzing Through Selective Symbolic Execution. In NDSS 2016 Feb (Vol. 16, pp. 1-16).
- [Sun2018]** Sun, Y., Wu, M., Ruan, W., Huang, X., Kwiatkowska, M., & Kroening, D. (2018). Concolic Testing for Deep Neural Networks. Proceedings of the 33rd ACM/IEEE International Conference on Automated Software Engineering, 109–119.
- [Taneja2008]** Taneja, K., & Xie, T. (2008). DiffGen: Automated Regression Unit-Test Generation. 2008 23rd IEEE/ACM International Conference on Automated Software Engineering, 407–410.
- [Tian2018]** Tian, Y., Pei, K., Jana, S., & Ray, B. (2018). DeepTest: Automated Testing of Deep-Neural-Network-Driven Autonomous Cars. Proceedings of the 40th International Conference on Software Engineering, 303–314.
- [Yang2012]** Yang, G., Pasareanu, C. S., & Khurshid, S. (2012). Memoized Symbolic Execution. Proceedings of the 2012 International Symposium on Software Testing and Analysis, 144–154.



**END OF DOCUMENT**