



# SoK: Prudent Evaluation Practices for Fuzzing

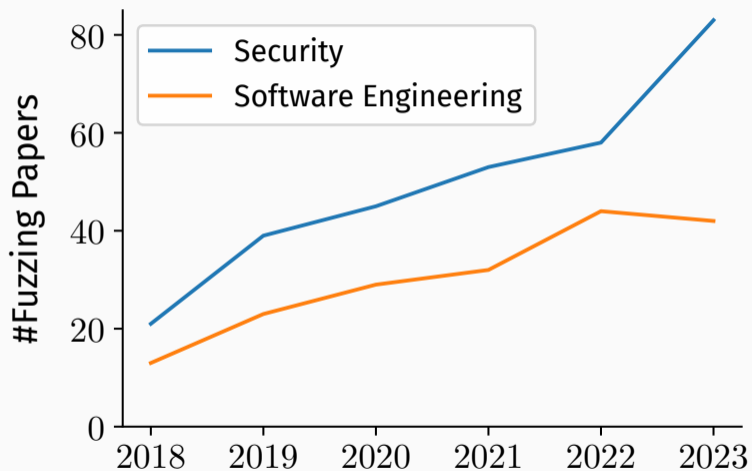
---

Moritz Schloegel, Nils Bars, Nico Schiller, Lukas Bernhard, Tobias Scharnowski, Addison Crump, Arash Ale Ebrahim, Nicolai Bissantz, Marius Muench, and Thorsten Holz

*CISPA Helmholtz Center for Information Security  
Ruhr University Bochum  
University of Birmingham*



## Fuzzing Papers are Still Popular



# Are we evaluating fuzzers right?

## SoK: Prudent Evaluation Practices for Fuzzing

Moritz Schloegel<sup>1</sup>, Nils Bars<sup>1</sup>, Nico Schiller<sup>1</sup>, Lukas Bernhard<sup>1</sup>, Tobias Scharnowski<sup>1</sup>  
Addison Crump<sup>1</sup>, Arash Ale-Ebrahim<sup>1</sup>, Nicolai Bissantz<sup>2</sup>, Marius Muench<sup>3</sup>, Thorsten Holz<sup>1</sup>

<sup>1</sup>*CISPA Helmholtz Center for Information Security, {first.lastname}@cispa.de*

<sup>2</sup>*Ruhr University Bochum, nicolai.bissantz@ruhr-uni-bochum.de*

<sup>3</sup>*University of Birmingham, m.muench@bham.ac.uk*

**Abstract**—Fuzzing has proven to be a highly effective approach to uncover software bugs over the past decade. After AFL popularized the groundbreaking concept of lightweight coverage

## 1. Introduction

*Fuzzing*, a portmanteau of “fuzz testing”, has gained

Study reproducibility

Study reproducibility

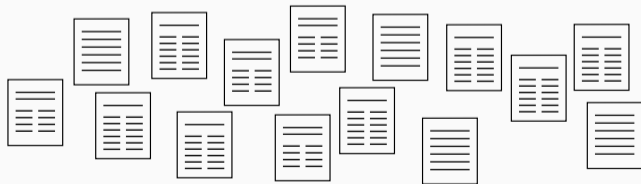
... and other evaluation pitfalls

# Approach

1. Find all fuzzing papers on 7 top-tier venues between 2018 and 2023

# Approach

1. Find all fuzzing papers on 7 top-tier venues between 2018 and 2023



# Approach

1. Find all fuzzing papers on 7 top-tier venues between 2018 and 2023
2. Literature analysis of 150 of these papers





# Approach

1. Find all fuzzing papers on 7 top-tier venues between 2018 and 2023
2. Literature analysis of 150 of these papers
3. Case studies: Reproducing experiments of 8 papers



# Approach

1. Find all fuzzing papers on 7 top-tier venues between 2018 and 2023
2. Literature analysis of 150 of these papers
3. Case studies: Reproducing experiments of 8 papers
4. Update recommendations (where needed)



# How to make sure fuzzing is reproducible?

Klees et al. – “Evaluating Fuzz Testing”, ACM CCS, 2018

Metzmann et al. – “FuzzBench: An Open Fuzzer Benchmarking Platform and Service”, ESEC/FSE, 2021

Böhme et al. – On the Reliability of Coverage Testing, ICSE, 2022

# How to make sure fuzzing is reproducible?

Klees et al. – “Evaluating Fuzz Testing”, ACM CCS, 2018

Metzmann et al. – “FuzzBench: An Open Fuzzer Benchmarking Platform and Service”, ESEC/FSE, 2021

Böhme et al. – On the Reliability of Coverage Testing, ICSE, 2022

No intention of finger pointing!

# Recommendations

- ① Document setup and parameters
- ② Sample relevant targets
- ③ Pick a good baseline
- ④ Choose suitable evaluation metrics
  - Code coverage
  - Bugs
- ⑤ Conduct a statistical evaluation

- 1 Document setup and parameters
- 2 Sample relevant targets
- 3 Pick a good baseline
- 4 Choose suitable evaluation metrics
  - Code coverage
  - Bugs
- 5 Conduct a statistical evaluation



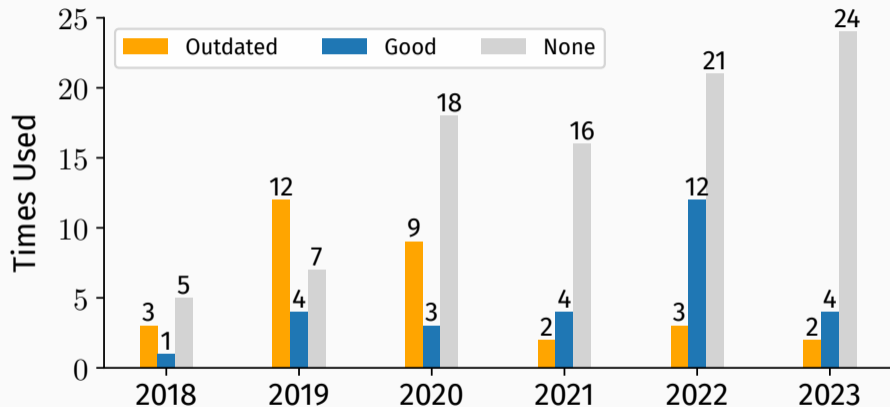
# Recommendations

- 1 Document setup and parameters
- 2 Sample relevant targets**
- 3 Pick a good baseline
- 4 Choose suitable evaluation metrics
  - Code coverage
  - Bugs
- 5 Conduct a statistical evaluation



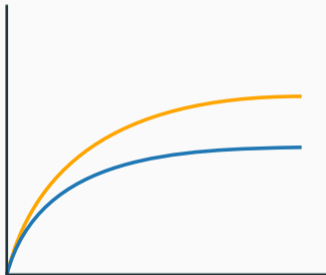


# Benchmarks



# Recommendations

- ① Document setup and parameters
- ② Sample relevant targets
- ③ Pick a good baseline**
- ④ Choose suitable evaluation metrics
  - Code coverage
  - Bugs
- ⑤ Conduct a statistical evaluation

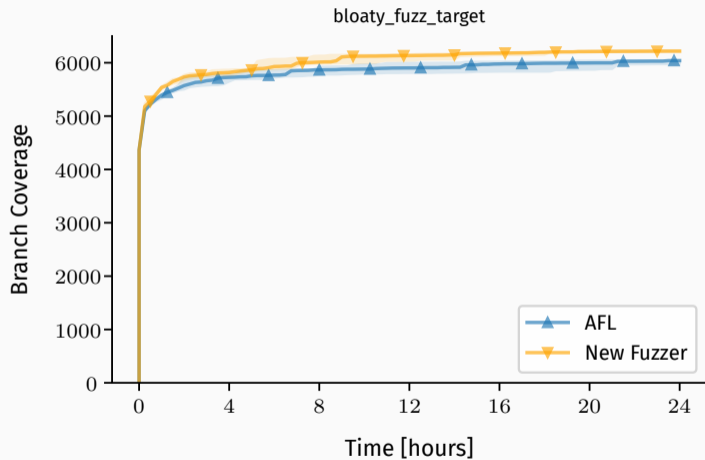


### ③ Pick a Good Baseline

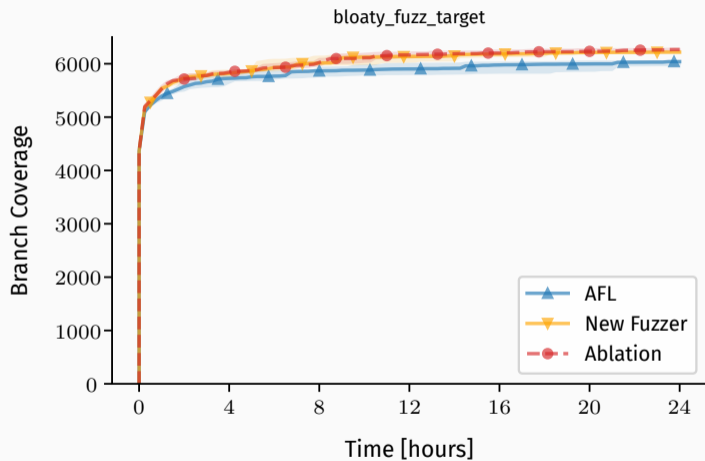
A new fuzzer that proposes to:

1. *Dynamically* adapt probabilities with which mutations are chosen
2. Use an *evolutionary strategy* to optimize these probabilities

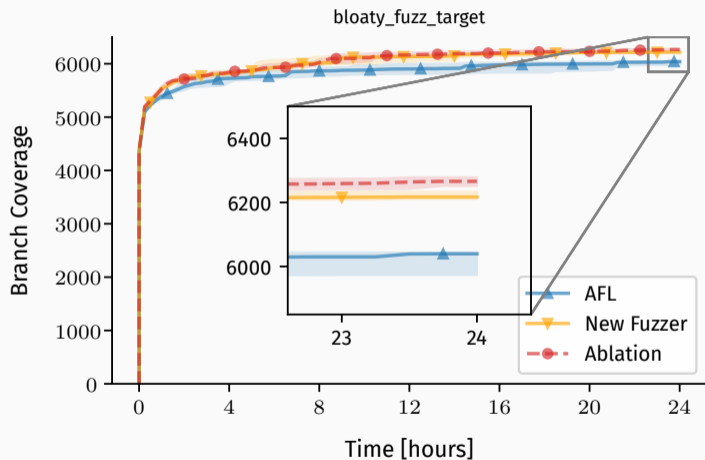
### ③ Pick a Good Baseline



### ③ Pick a Good Baseline



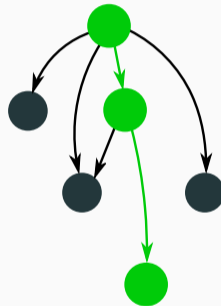
### ③ Pick a Good Baseline



Lesson learned: ablation studies are important

# Recommendations

- 1 Document setup and parameters
- 2 Sample relevant targets
- 3 Pick a good baseline
- 4 Choose suitable evaluation metrics**
  - Code coverage
  - Bugs
- 5 Conduct a statistical evaluation

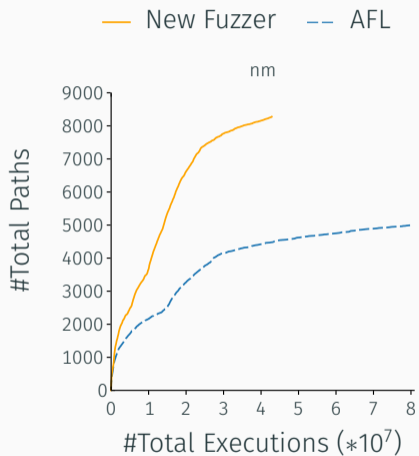




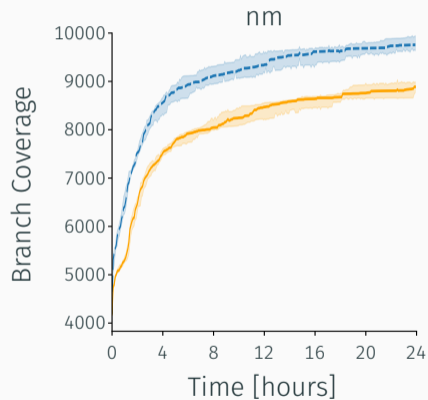
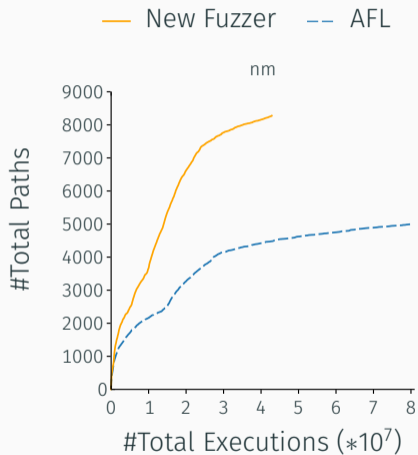
New fuzzer:

- Improves input scheduling
- Aims to covers more *paths* with fewer inputs

## ④ Suitable Metrics: Code Coverage



## ④ Suitable Metrics: Code Coverage



Lesson learned: new metrics may mislead readers

⇒ include known metrics!

# Recommendations

- 1 Document setup and parameters
- 2 Sample relevant targets
- 3 Pick a good baseline
- 4 Choose suitable evaluation metrics**
  - Code coverage
  - **Bugs**
- 5 Conduct a statistical evaluation

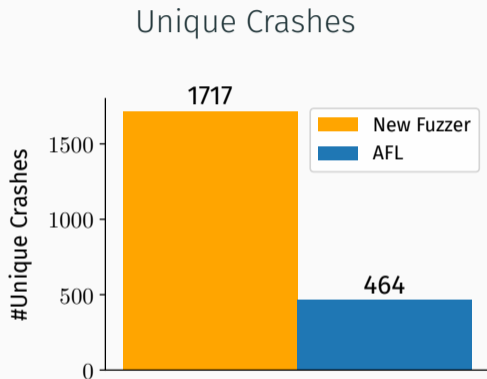


## ④ Suitable Metrics: Unique Crashes

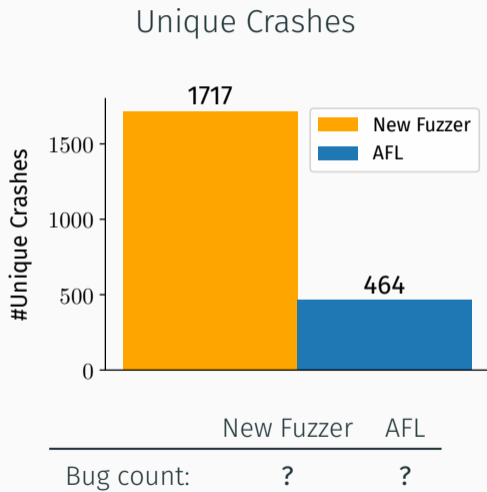
A new fuzzer that:

- Proposes memory usage as additional feedback
- Uses *unique crashes* as a metric

## ④ Suitable Metrics: Unique Crashes

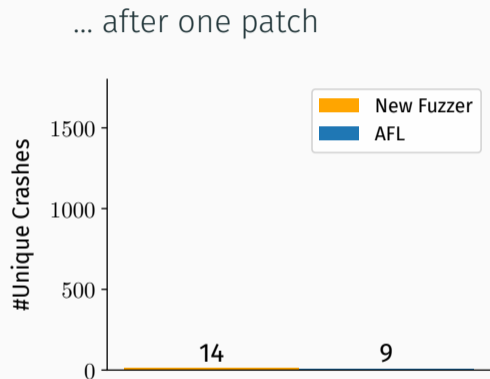
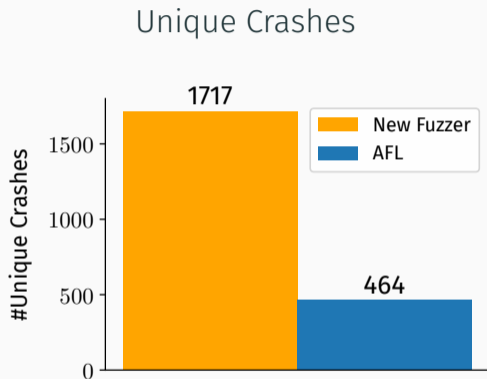


## ④ Suitable Metrics: Unique Crashes

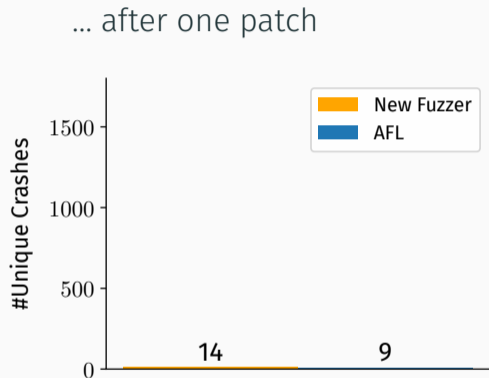
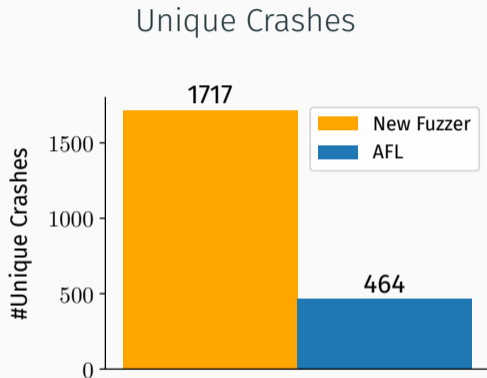




## ④ Suitable Metrics: Unique Crashes

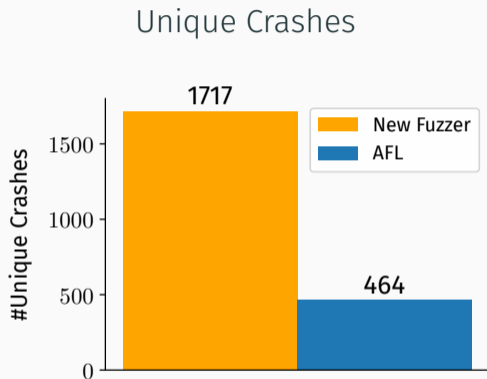


## ④ Suitable Metrics: Unique Crashes

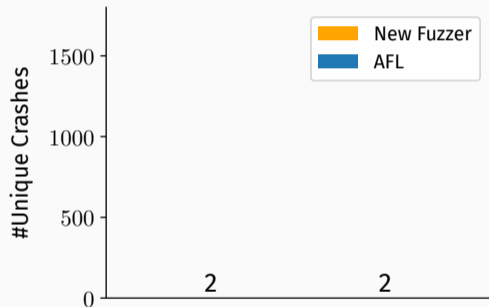


	New Fuzzer	AFL
Bug count:	1+?	1+?

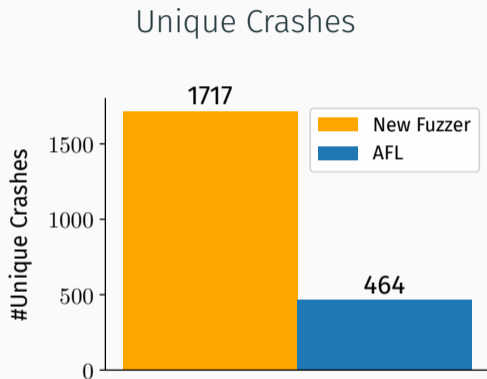
## ④ Suitable Metrics: Unique Crashes



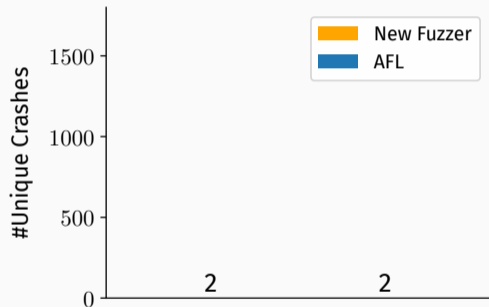
...after manual deduplication



## ④ Suitable Metrics: Unique Crashes



...after manual deduplication



	New Fuzzer	AFL
Bug count:	2	2

Lesson learned: “unique” crashes  $\neq$  actual bugs

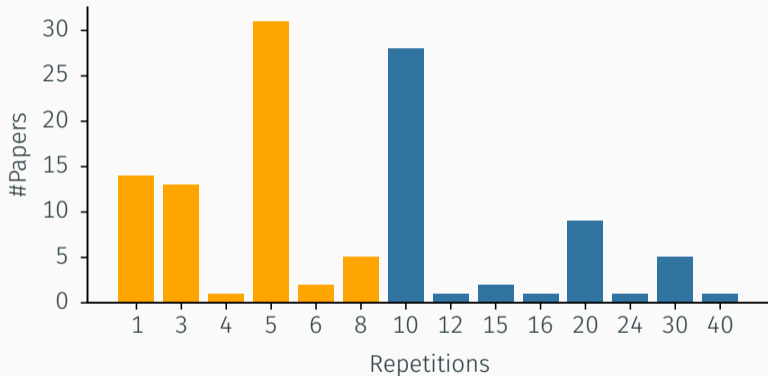
$\Rightarrow$  need deduplication or should use actual bugs

# Recommendations

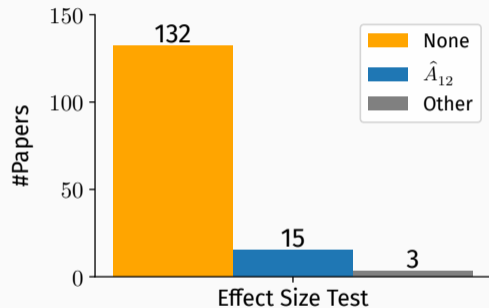
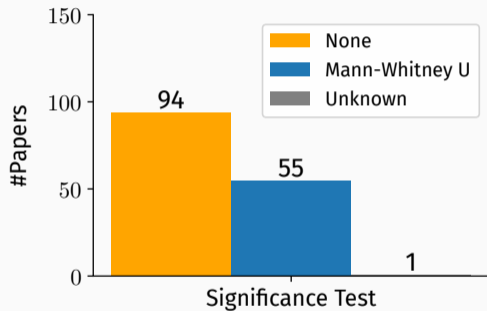
- 1 Document setup and parameters
- 2 Sample relevant targets
- 3 Pick a good baseline
- 4 Choose suitable evaluation metrics
  - Code coverage
  - Bugs
- 5 Conduct a statistical evaluation**



## ⑤ Statistically evaluate results



## ⑤ Statistically evaluate results



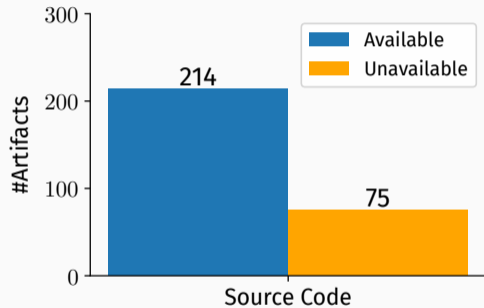


There's more beyond the evaluation itself:

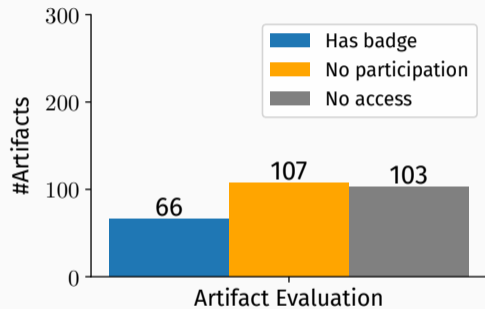
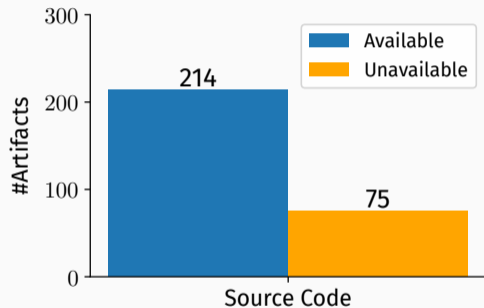
- What about the fuzzer source code?
- What about new bugs found during the evaluation?

Beyond the paper: artifact availability

# Artifact Availability



# Artifact Availability



Good: much code is openly available!

BUT: low artifact evaluation participation

There's more beyond the evaluation itself:

- What about the fuzzer source code?
- What about new bugs found during the evaluation?

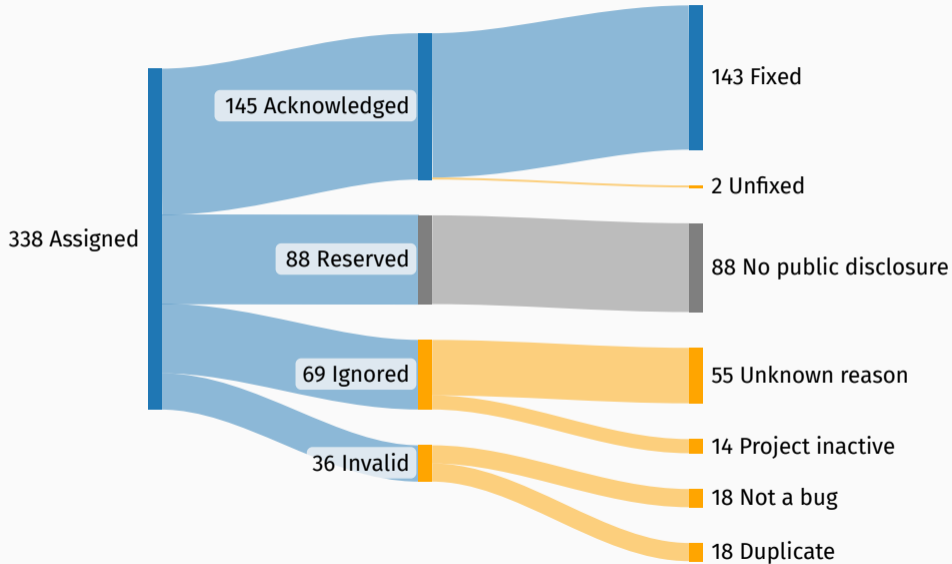
What happens to found bugs?

What happens to found bugs?

⇒ Responsible disclosure?



- ① Look for CVEs in fuzzing papers
- ② Check their outcome



Misaligned incentives

Misaligned incentives & no verification

Misaligned incentives & no verification

⇒ Easy to game the system

# Recommendations

- 1 Document setup and parameters
  - 2 Sample relevant targets
  - 3 Pick a good baseline
  - 4 Choose suitable evaluation metrics
    - Code coverage
    - Bugs
  - 5 Conduct a statistical evaluation
- + Artifact availability
- + CVE misuse

} ⇒ guidelines

## Fuzzing Evaluation Guidelines

Current version: 1.0.3

Proposals for changes welcome (please open an issue for discussion or a pull request for changes).

DISCLAIMER: These items represent are a best-effort attempt at capturing action items to follow during the evaluation of a scientific paper that focuses on fuzzing. **They do not apply universally to all fuzzing methods - in certain scenarios, techniques may wish to deviate for good reason from these guidelines. In any case, a case-by-case judgment is necessary.** The guidelines do not discuss many malicious choices that immediately negate any chance of a fair evaluation, such as giving your fuzzer an unfair advantage (e.g., by fine-tuning the fuzzer or its targets) or putting other fuzzers at a disadvantage.

### A. Preparation for Evaluation

1. Find relevant tools and baselines to compare against
  - 1.1 Include state-of-the-art techniques from both academia and industry
  - 1.2 If your fuzzer is based on an existing fuzzer, include the baseline (to measure the delta of your changes, which allows attributing improvements to your technique)
  - 1.3 Use recent versions of fuzzers
  - 1.4 If applicable, derive a baseline variant of your technique that replaces core contributions by alternatives. For example, consider using a variant that replaces an informed algorithm with randomness.

① Fuzzing evaluations are hard to get right

② Join artifact evaluation

③ Help us shape the guidelines



Paper



Guidelines