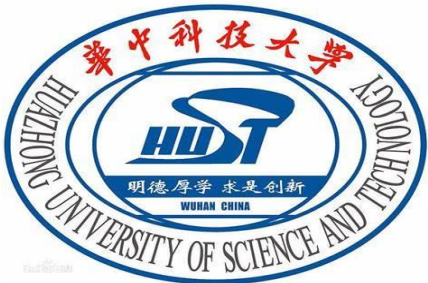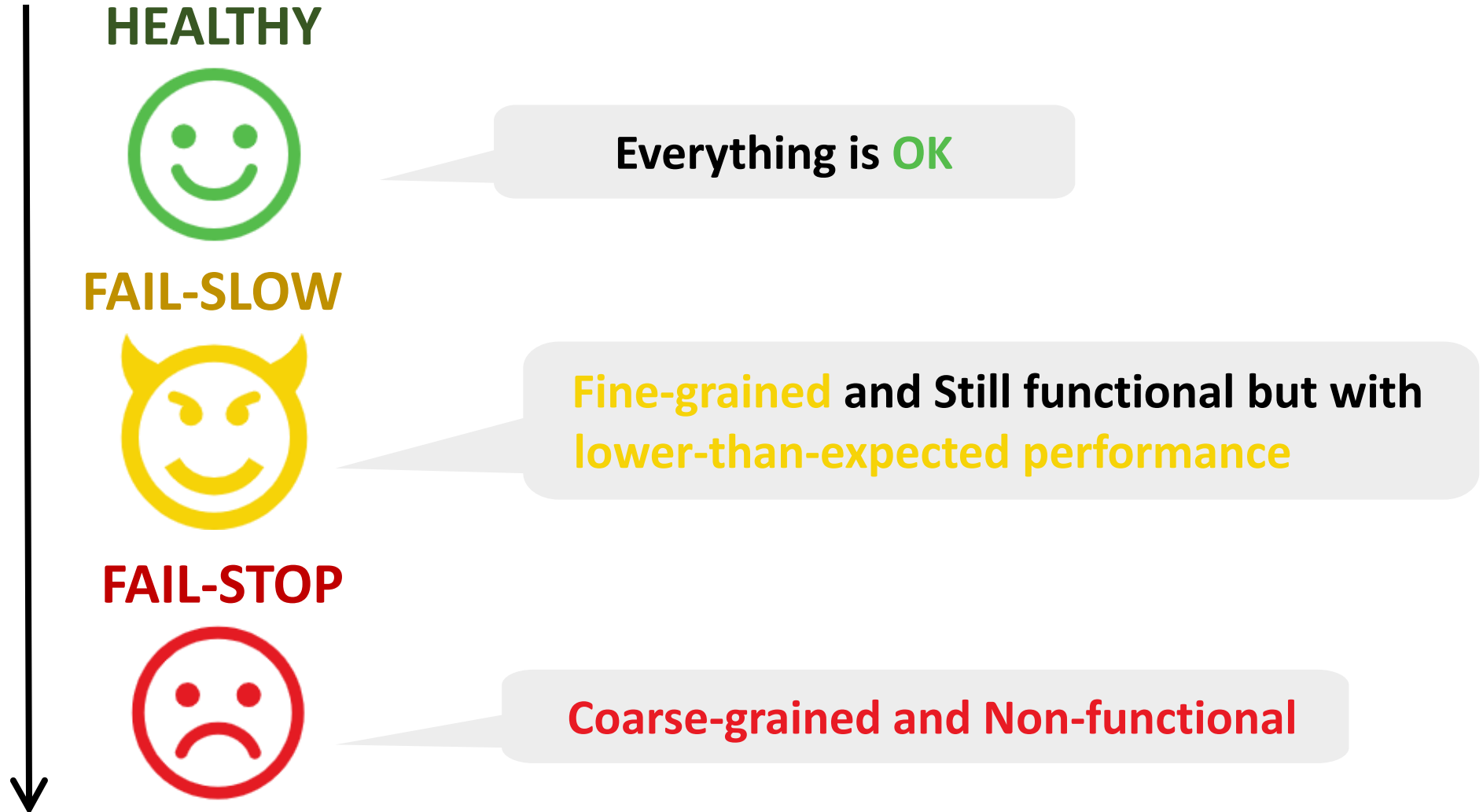# Understanding and Detecting Fail-Slow Hardware Failure Bugs in Cloud Systems

**Gen Dong**, Yu Hua, Yongle Zhang*, Zhangyu Chen, Menglei Chen
Huazhong University of Science and Technology
*Purdue Univeristy

USENIX ATC 2025

# Hardware Failures in the Wild

**HEALTHY**

Everything is OK

**FAIL-SLOW**

Fine-grained and Still functional but with lower-than-expected performance

**FAIL-STOP**

Coarse-grained and Non-functional

# Fail-Slow Hardware is a Real-World Problem

**FAIL-SLOW**

**Severe**

A 1Gb NIC card on a machine that suddenly only transmits at 1 kbps[1]

Fail-slow NVMe SSDs can degrade to SATA SSD or HDD-level performance[2]
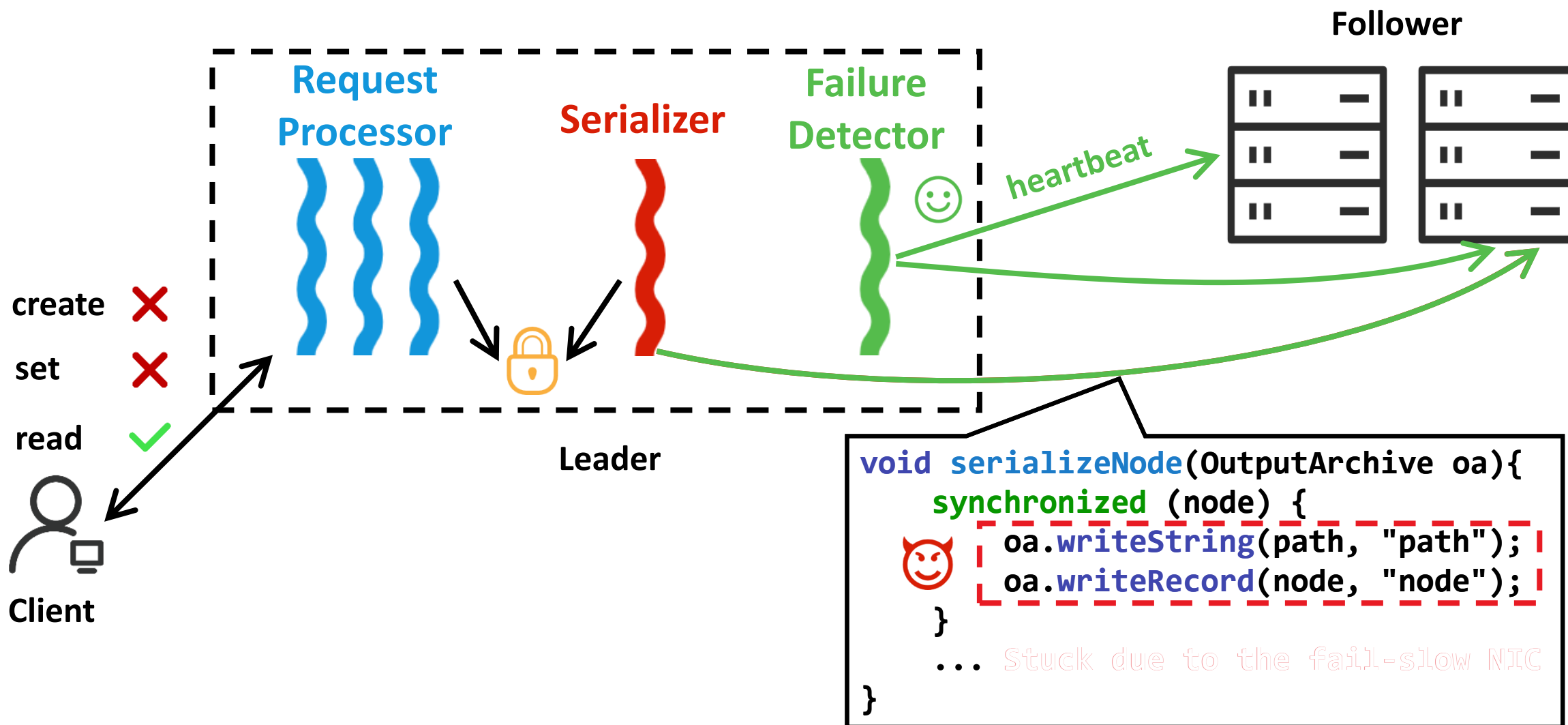
As frequent as fail-stop incidents[2]

**Common**

Annual fail-slow failure rate is 1-2%[3]

[1] Fail-Slow at Scale: Evidence of Hardware Performance Faults in Large Production Systems, Guanwai et al.
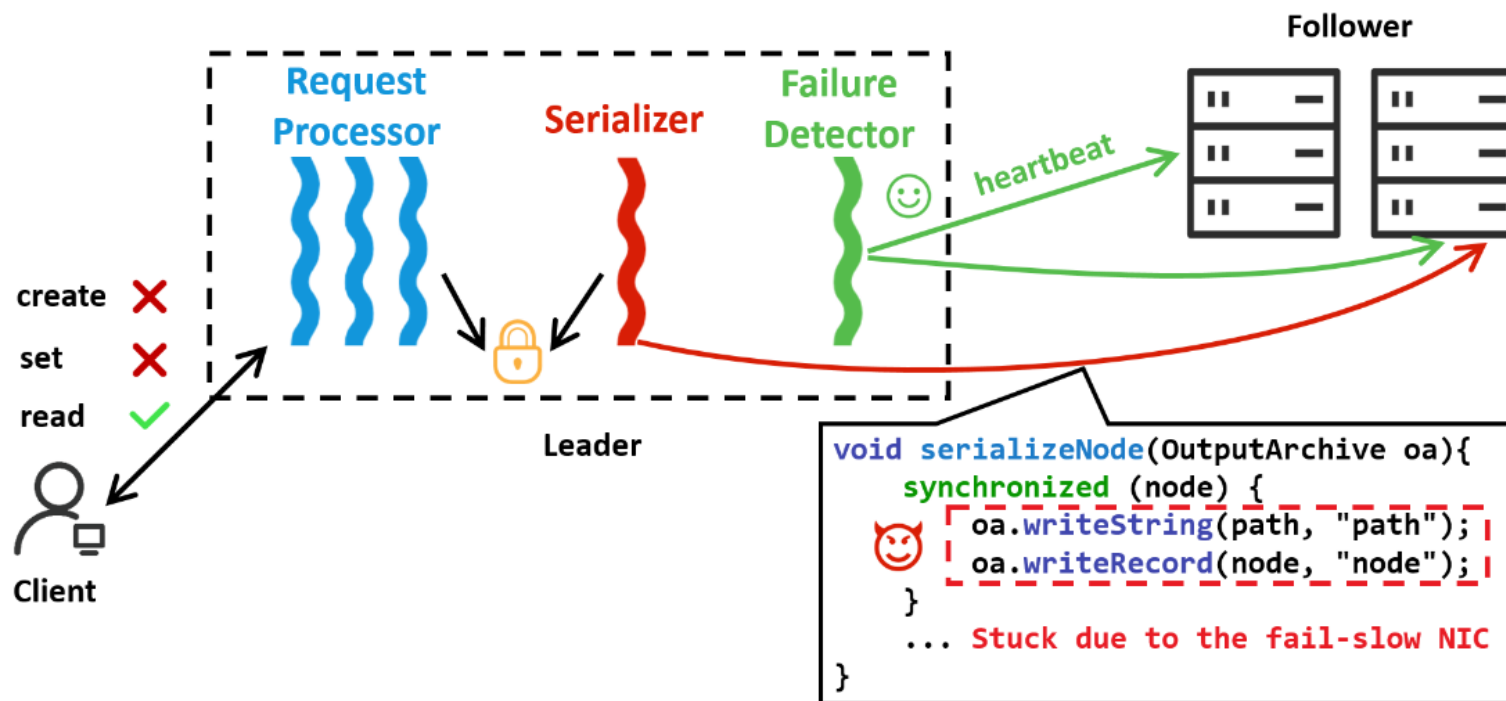[2] NVMe SSD Failures in the Field: the Fail-Stop and the Fail-Slow, Lu et al.
[3] IASO: A Fail-Slow Detection and Mitigation Framework for Distributed Storage Services, Panda et al.

# A Real Bug in ZooKeeper



**Follower**

Request Processor

Serializer

Failure Detector ☺

heartbeat

create ✗
set ✗
read ✓

Client

Leader

```
void serializeNode(OutputArchive oa){
    synchronized (node) {
        oa.writeString(path, "path");
        oa.writeRecord(node, "node");
    }
    ...
}
```
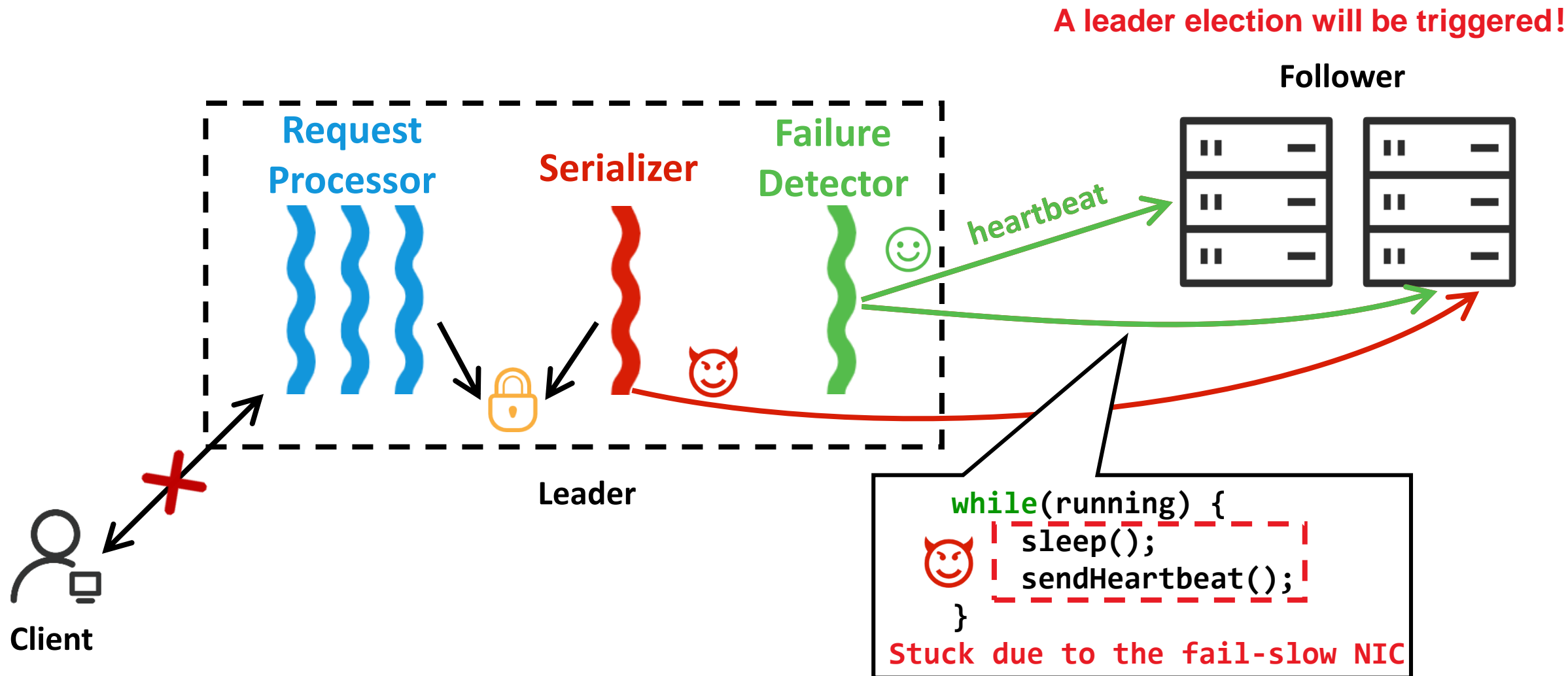
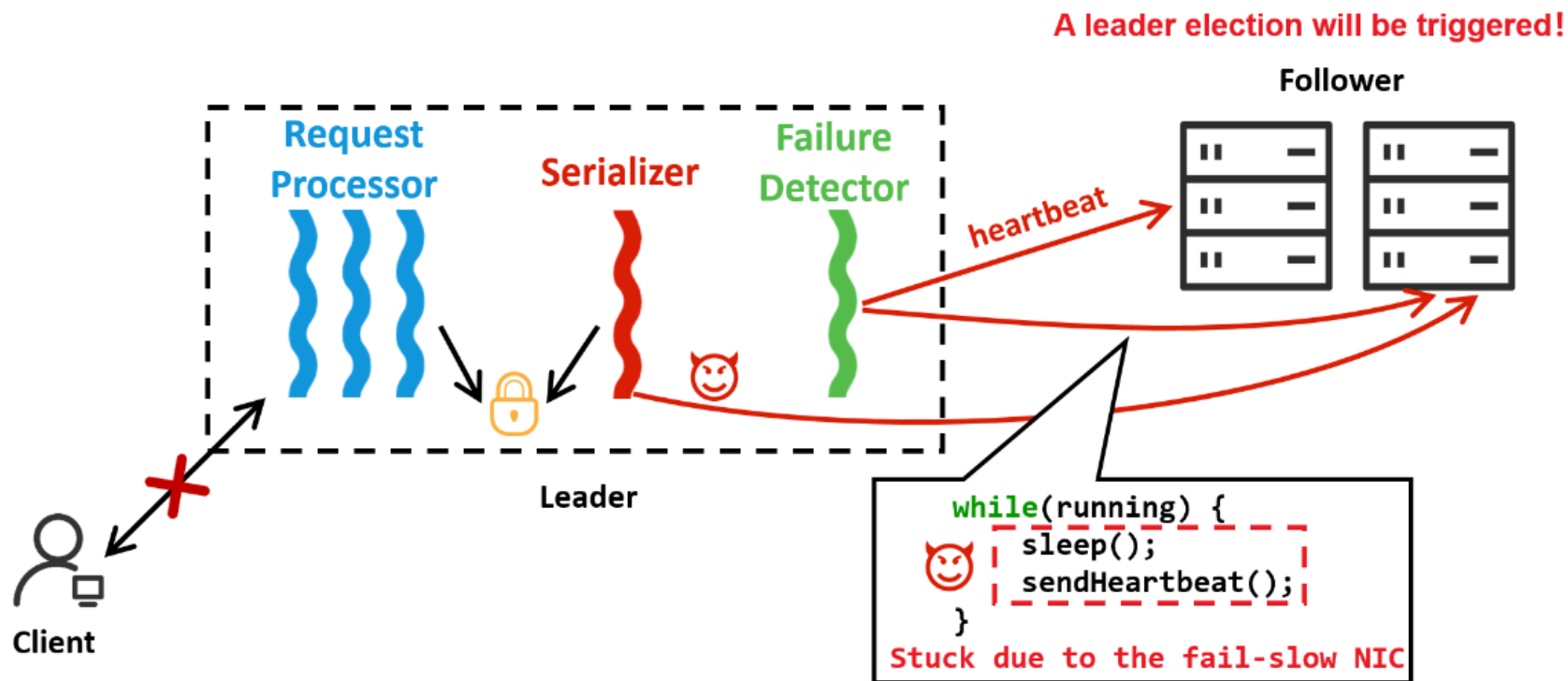Stuck due to the fail-slow NIC

4

# A Real Bug in ZooKeeper



We define **fail-slow hardware failures (FSH failure)** as software-level failures caused by fail-slow hardware.

# A Real Bug in ZooKeeper



A leader election will be triggered!

Follower

Request Processor

Serializer

Failure Detector

heartbeat

Leader

Client

```
while(running) {
    sleep();
    sendHeartbeat();
}
```
Stuck due to the fail-slow NIC

6

# A Real Bug in ZooKeeper



The fine granularity of fail-slow hardware is necessary to trigger FSH failures **(a subset of I/O operations)**
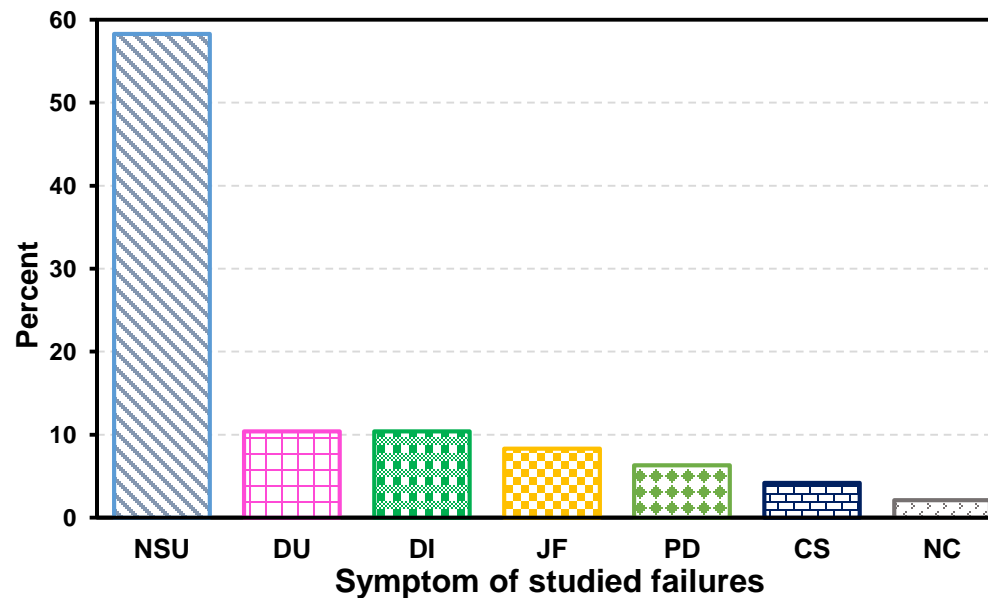
# Study Methodology

- We study 48 FSH failure cases from five large, widely-used cloud systems.
  - Diverse services
    - Coordination service, file system, data-analytic framework, and database

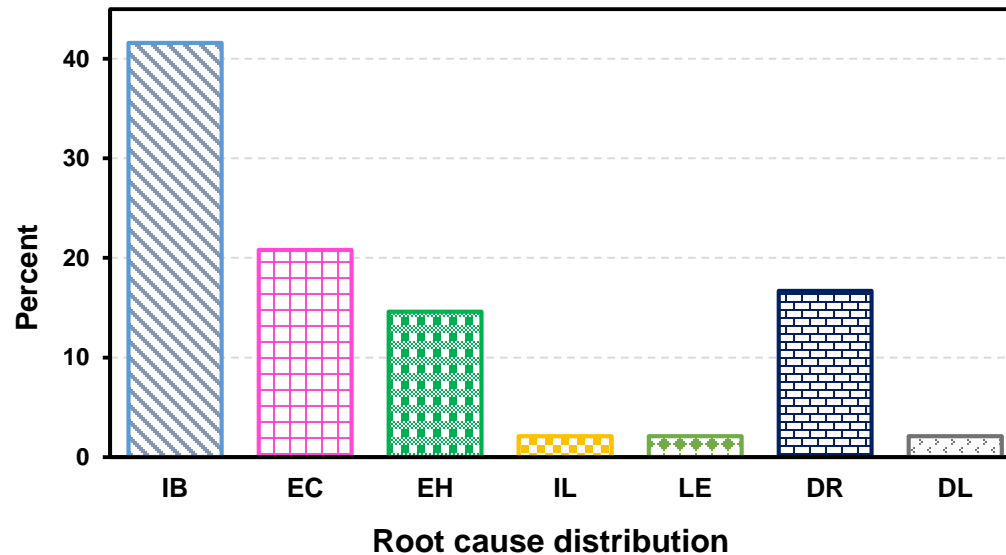| Systems | Cases | Versions | Date |
|---|---|---|---|
| ZooKeeper | 11 | 16 | 2009/05/27-2023/10/13 |
| HDFS | 18 | 25 | 2012/07/02-2022/09/07 |
| HBase | 10 | 18 | 2014/03/24-2023/12/16 |
| MapReduce | 4 | 3 | 2010/05/20-2022/05/22 |
| Cassandra | 5 | 7 | 2010/08/26-2020/12/09 |

# Understanding FSH failures

- Finding 1: over half (58.3%) of FSH failures cause node service to be unavailable.

- Finding 2: 20.8% of FSH failures are silent (including data unavailability and inconsistency).

NSU: node service unavailable; DU: data unavailability; DI: data inconsistency; JF: job failure
PD: performance degradation; CS: client stuck; NC: node crash
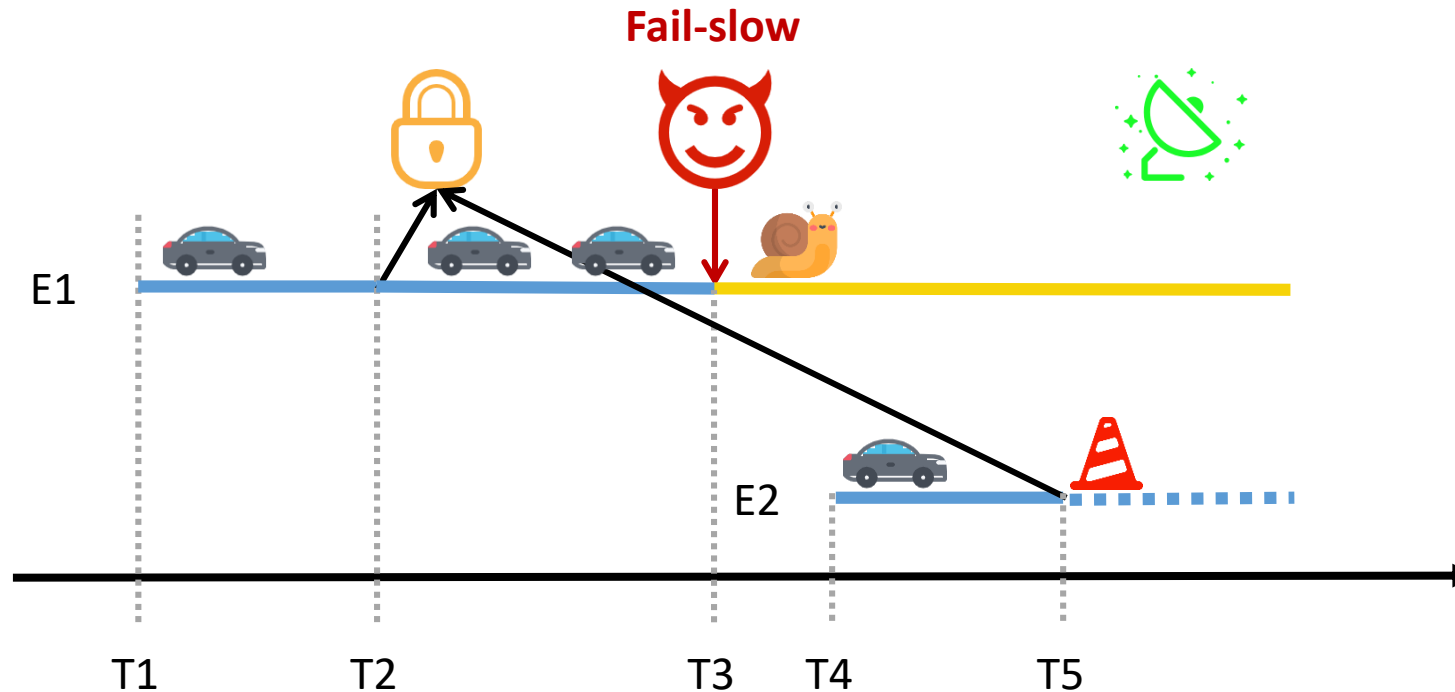
9

# Understanding FSH failures

- Finding 3
  - Root causes are diverse.
  - The top three (total 93.7%) root causes are indefinite blocking, buggy internal checker, and data race.



**Synchronized and timeout mechanisms are vulnerable.**

# Understanding FSH failures

- Synchronized mechanisms are vulnerable
- Fine granularity of fail-slow hardware is necessary

# How to Deal with FSH failures

- Existing in-production detectors
  - Panorama[OSDI'18]
  - IASO[ATC'19]
  - OmegaGen[NSDI'20]
  - PERSEUS[FAST'23]

  FSH failures already cause damages!

- Existing fault injection tools
  - FATE[NSDI'11]
  - CrashTuner[SOSP'19]
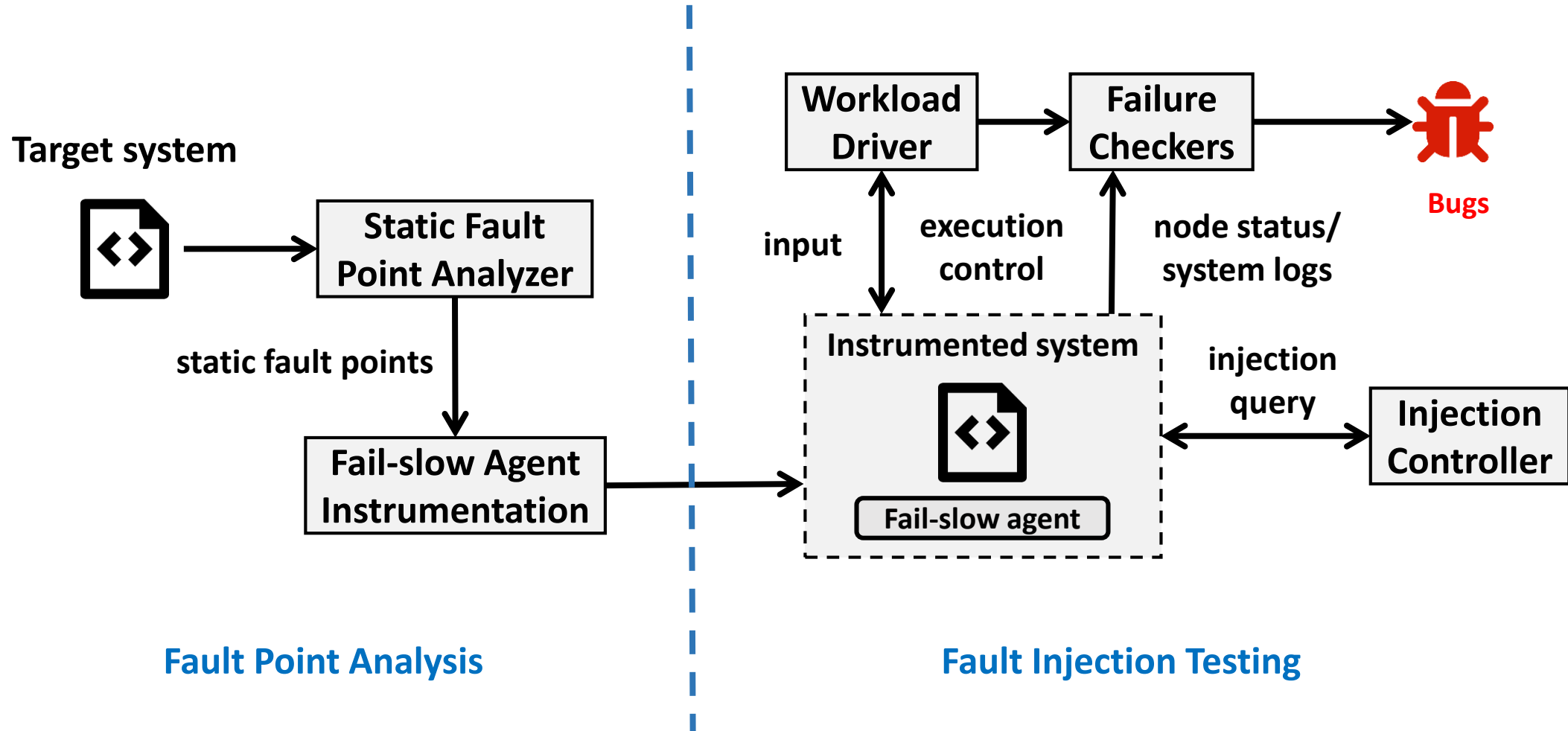  - Legolas[NSDI'24]
  - Chronos[S&P'24]

  Overlooking characteristics of FSH failures!

# Our Solution: Sieve

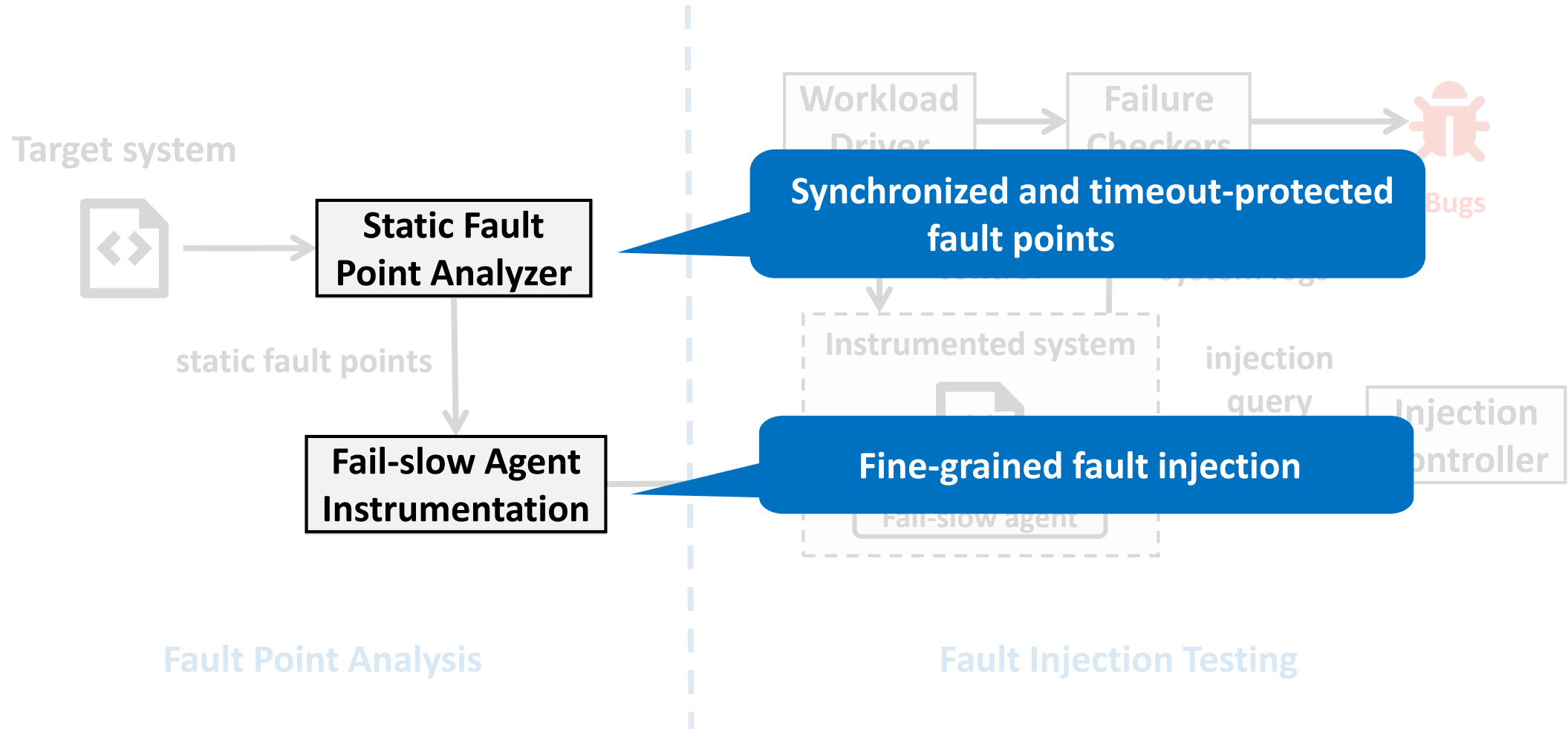> **A fault injection testing framework for cloud systems to detect FSH failures**

- Efficiently explore the large fault injection space
  - Statically analyze synchronized and timeout-protected fault points
- Enable fine-grained fault injection
  - Automatically instrument hooks to precisely simulate fail-slow hardware within a system

# Sieve Workflow



**Fault Point Analysis**

**Fault Injection Testing**

# Sieve Workflow



Target system

**Static Fault Point Analyzer**

static fault points

**Fail-slow Agent Instrumentation**

Workload Driver

Failure Checkers

**Synchronized and timeout-protected fault points**

Bugs

Instrumented system

injection query

Injection Controller

**Fine-grained fault injection**

Fail-slow agent

Fault Point Analysis

Fault Injection Testing

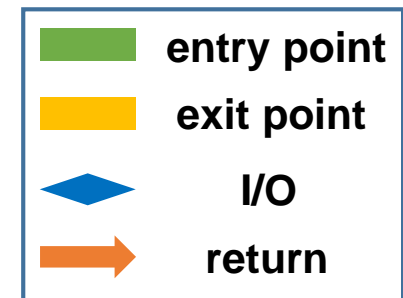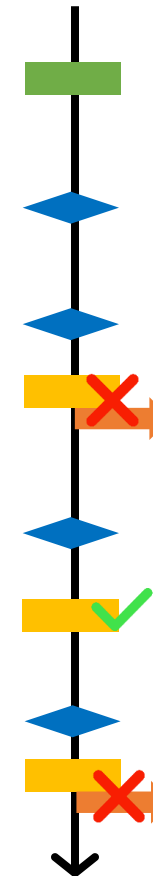# Static Fault Point Analyzer

- Identify synchronized fault points

```
1   if(cond1){
2       synchronized(…){
3               if(cond2){
4                       I/O₁;
5               }else{
6                       I/O₂;
7                       return …;
8               }
9               I/O₃;
10      }                    Critical region
11  }else{…;}
12  I/O₄;
13  return …;
```
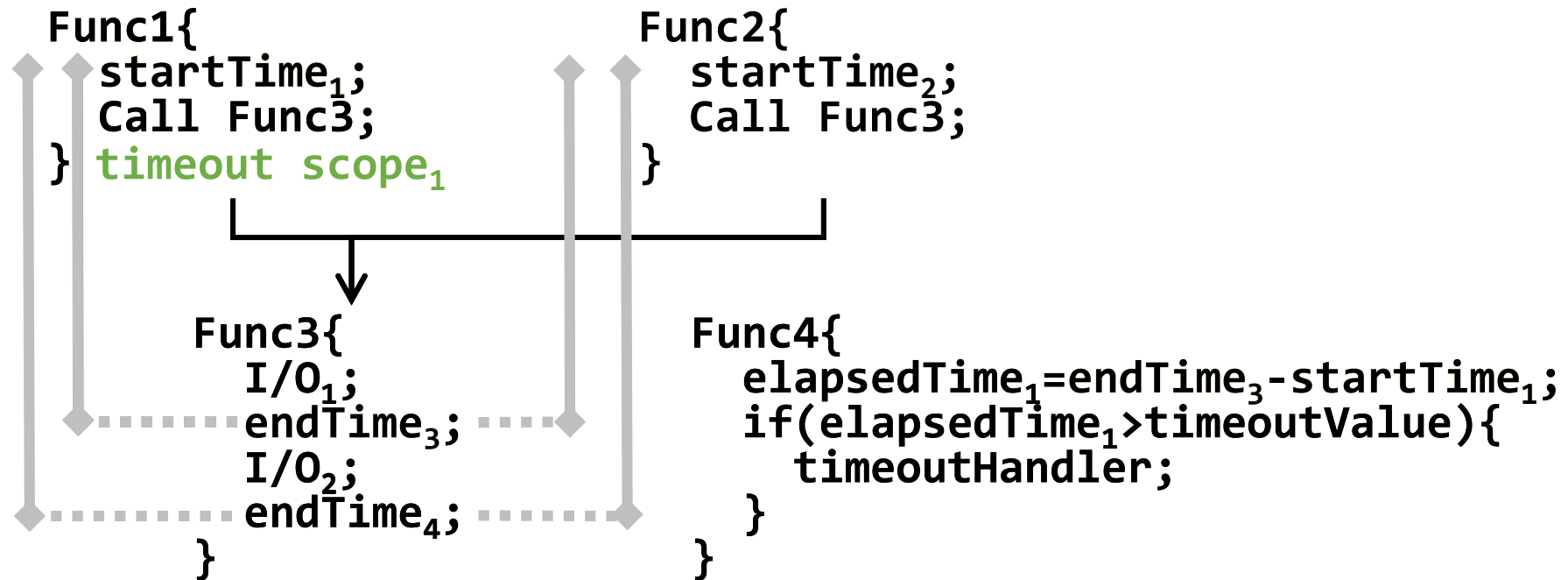


translate

entry point
exit point
I/O
return

1. the last exit point
2. not followed by return

16

# Static Fault Point Analyzer

- Identify timeout-protected fault points

```
Func1{                          Func2{
  startTime₁;                     startTime₂;
  Call Func3;                     Call Func3;
} timeout scope₁                }



        Func3{                    Func4{
          I/O₁;                     elapsedTime₁=endTime₃-startTime₁;
          endTime₃;                 if(elapsedTime₁>timeoutValue){
          I/O₂;                       timeoutHandler;
          endTime₄;                 }
        }                         }
```

# Fail-Slow Agent Instrumentation

- Coarse-grained vs. Fine-grained



**Node level**

**API level**

**Implementation level**

Controlled env

Interception code

Library

◆ fail-slow agent

Normal env and library
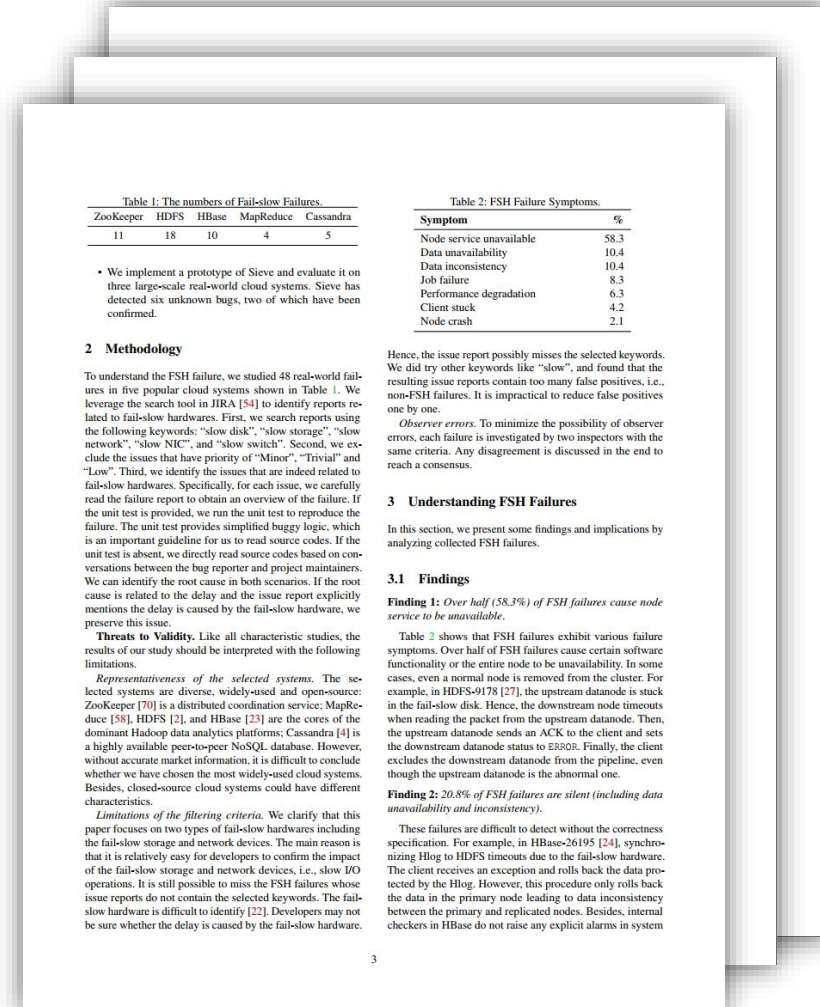
Injection Controller

Sieve

# Evaluation

- Applied Sieve to three cloud systems
  - ZooKeeper, HDFS, Kafka

- Can Sieve effectively find new bugs?
  - Detected six new bugs, two of which are confirmed

| Bug ID | Failure Symptoms | Status |
|--------|------------------|--------|
| ZK-4816 | A follower cannot follow the leader for a long time | Pending |
| ZK-4817 | CancelledKeyException cannot catch the client disconnection exception | Pending |
| ZK-4844 | Fail-slow disk while executing writeLongToFile causes the follower to hang | Pending |
| ZK-4836 | Inconsistent ACL index leads to MarshallingError | Confirmed |
| KA-16401 | One request consumes all request handler threads | Pending |
| KA-16412 | An uncreated topic is considered as a created one | Confirmed |

# More Details

- More bug study details
- Fault injection strategies
- Bug explanation
- ......

# Conclusion

- Fail-slow hardware causes severe damages in cloud systems
  - Existing fault injection testing is inefficient
- We conduct a study on 48 FSH failure cases
- Sieve: a fault injection testing framework to detect FSH failure bugs
  - Identify synchronized and timeout-protected fault points
  - Enable fine-grained fault injection
- Found six bugs, two of which are confirmed
- Open Sourced at https://github.com/RabbitDong-on

# Thank you! Q&A