

# Monarch: A Fuzzing Framework for Distributed File Systems

Tao Lyu Liyi Zhang Zhiyao Feng Yueyang Pan Yujie Ren  
Meng Xu Mathias Payer Sanidhya Kashyap

**EPFL**



UNIVERSITY OF  
**WATERLOO**

# DFSes are critical infrastructure



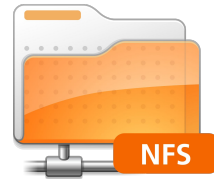
High Performance  
Computing



Serverless<sup>1</sup>



Machine Learning  
Model Training<sup>2</sup>



...

# But DFSes are not reliable

**11K Ceph bug reports**

**17K Lustre bug reports**

#	Project	Tracker	Status	Priority	Subject	Assignee	Updated
66	Lustre						dvanced
66		Order by					
66			LU-17917		sanity-hsm: 26A ([: != 5 : syntax error: operand expected (error token is "!= 5 ")		
66			LU-17916		Enable folio allocation support on the buffed io read/write path (BIO)		
66			LU-17915		gds io crashes in unaligned_dio		
66			LU-17914		Inetctl net set command issues false error		
66			LU-17913		sanity-lnet test_220 is silently failing		
66			LU-17911		Faked flexible array usage causes crash when Fortify feature is enabled		

Memory-unsafe languages

State-sharing across clients

High concurrency

Fault tolerance

Finding and fixing bugs in DFSes become important

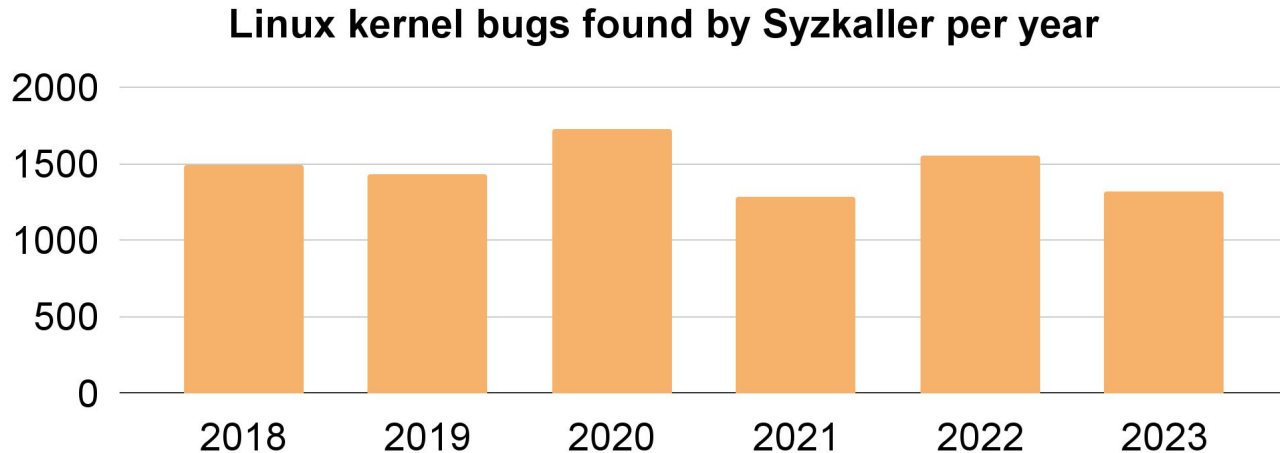
# Existing approaches to finding DFSEs bugs

<b>Regression Testing</b>	<b>Model Checking</b>	<b>Formal Verification</b>
Test suit	Modist SAMC	Verdi Ironfleet
<b>Manual expert effort</b> <b>Limited test cases</b>	<b>State explosion</b>	<b>Manual expert effort</b> <b>Mostly verify critical parts</b>

**Any automatic and scalable bug-finding techniques?**

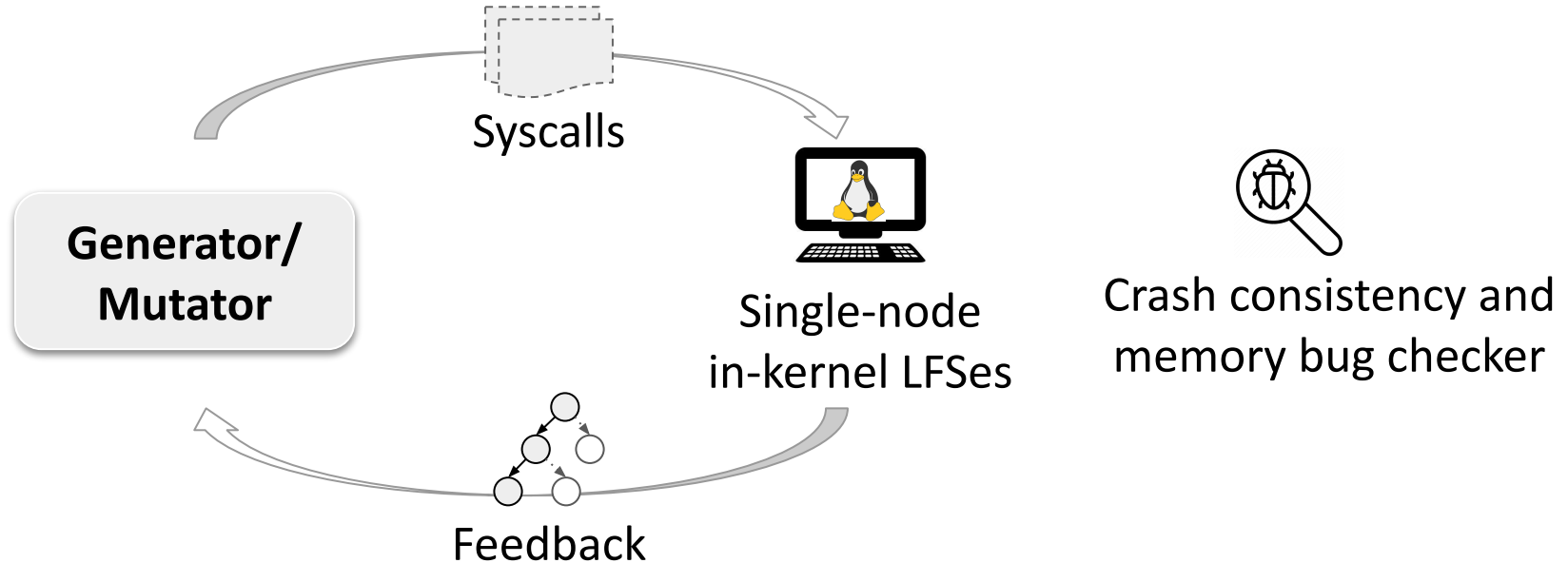
# Fuzzing: A practical and impactful approach

- Linux kernel fuzzer, **syzkaller**, detects over **1K bugs every year**<sup>1</sup>
- Specifically, more than **800 bugs in local file systems** over the years

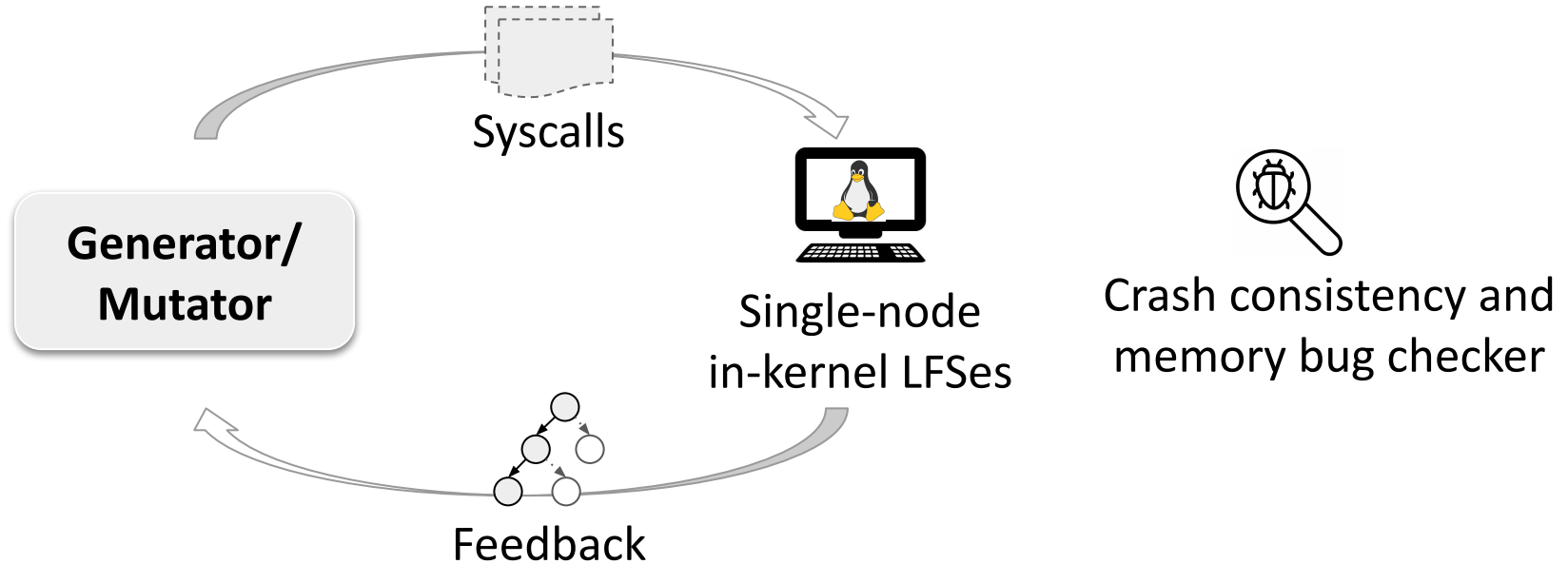


1. Syzbot/Syzkaller <https://syzkaller.appspot.com/upstream/graph/found-bugs>

# How does fuzzing work?



# How does fuzzing work?



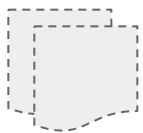
**Current LFS fuzzers are not applicable for DFSes**

# Missing pieces for fuzzing DFSes



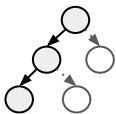
Single-node and  
in-kernel file system

**Multi-node cross-kernel/userspace fuzzing  
architecture**



Syscalls

**Distributed faults as a testing input space**



Representation of  
single-kernel exe state

**Representation of cross-node and  
cross-kernel/userspace execution states**



Crash consistency +  
memory checker

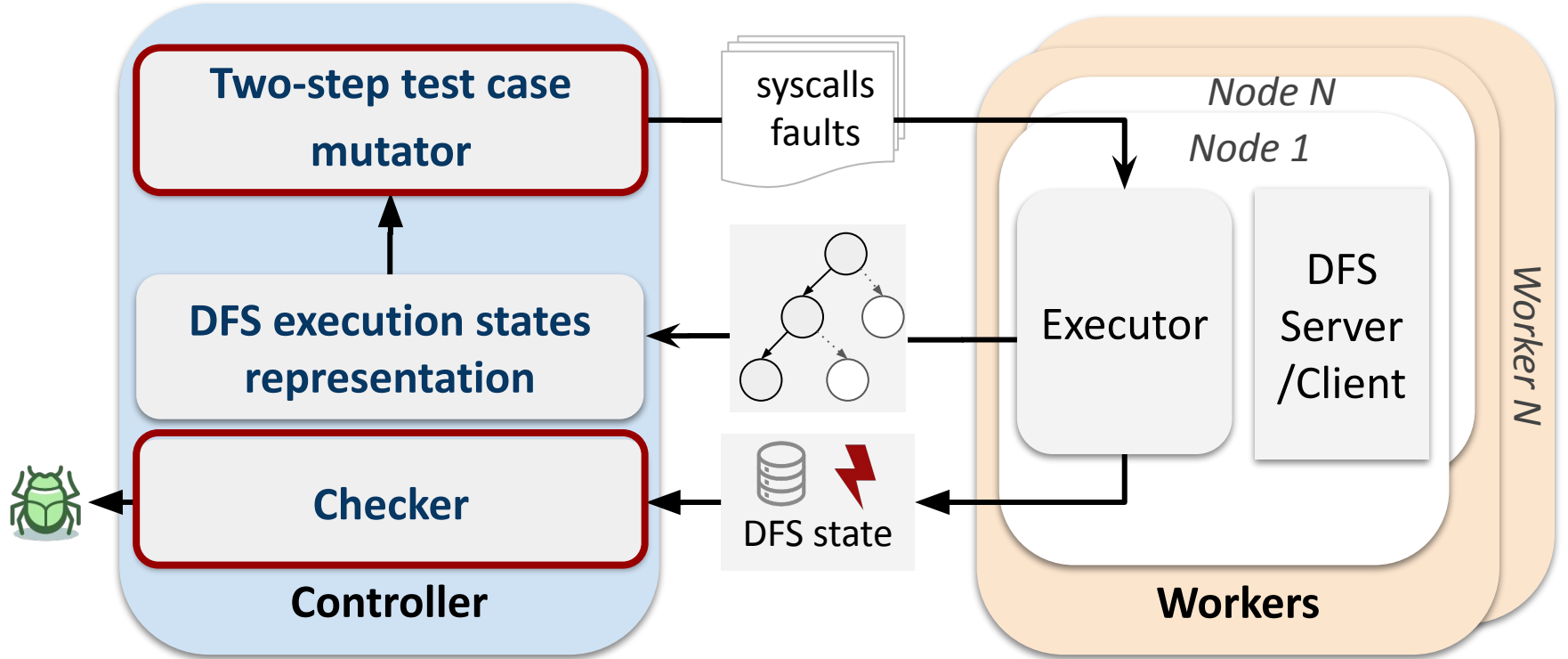
**A systematic DFS semantic checker**



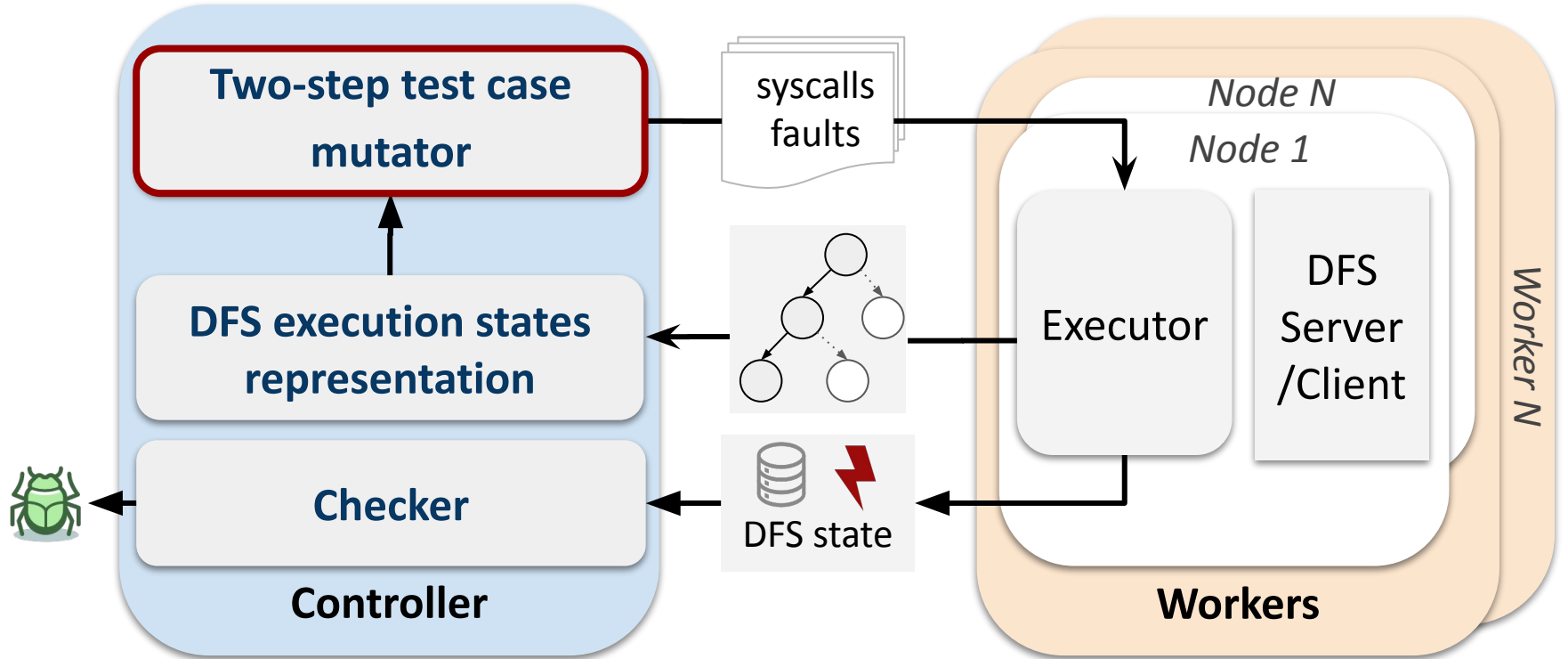
# **MONARCH**

## **A distributed file system fuzzer**

# Monarch architecture



# Monarch architecture



# Two-step test case mutator

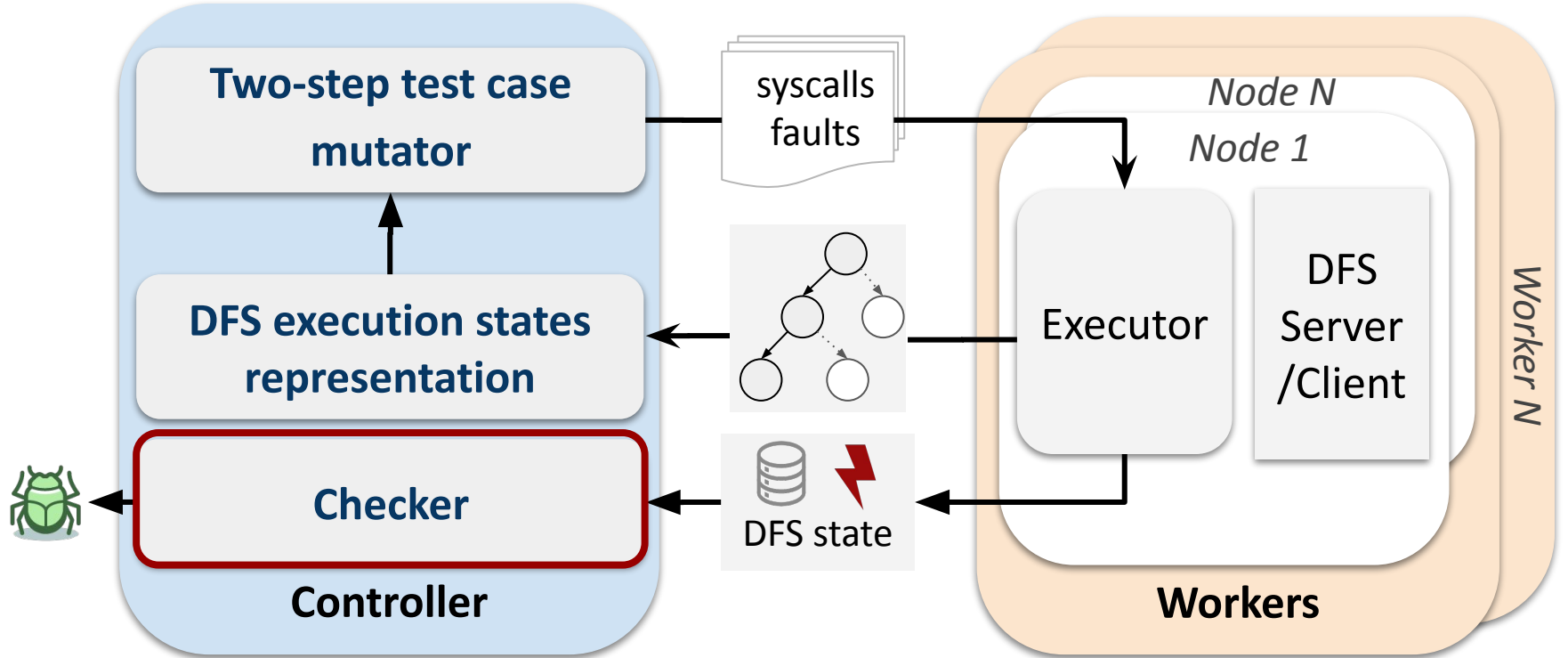
Non-fault mode: Testing with file system syscalls only (e.g., open, read)

- Syscall templates and mutation rules

Fault mode: Testing with syscalls + injected faults

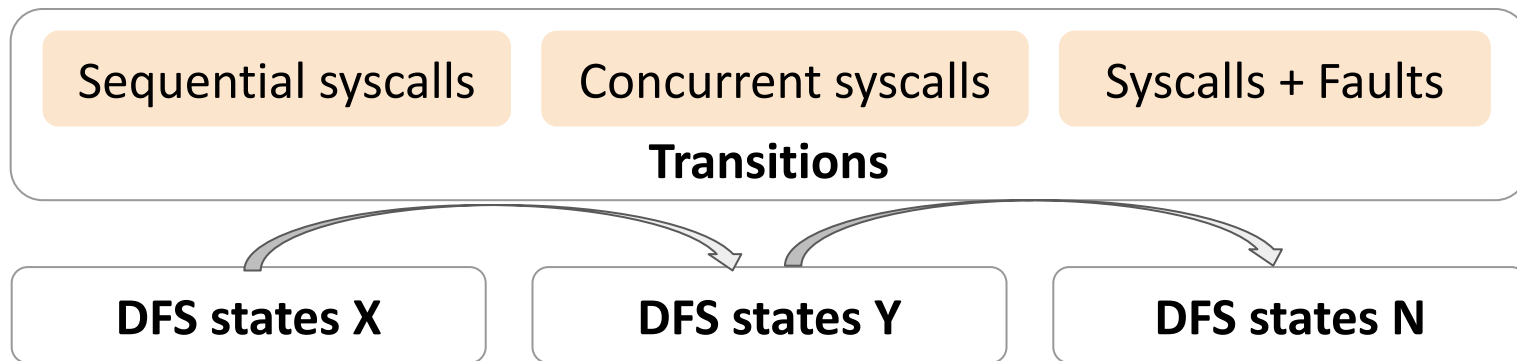
- **Distributed faults** → Network partitions, node crashes
- **How to deterministically inject faults?** → Synchronization primitives
- **At which granularity to inject faults?** → Syscall granularity
- **When to inject faults?** → After a test case triggers new execution state

# Monarch architecture



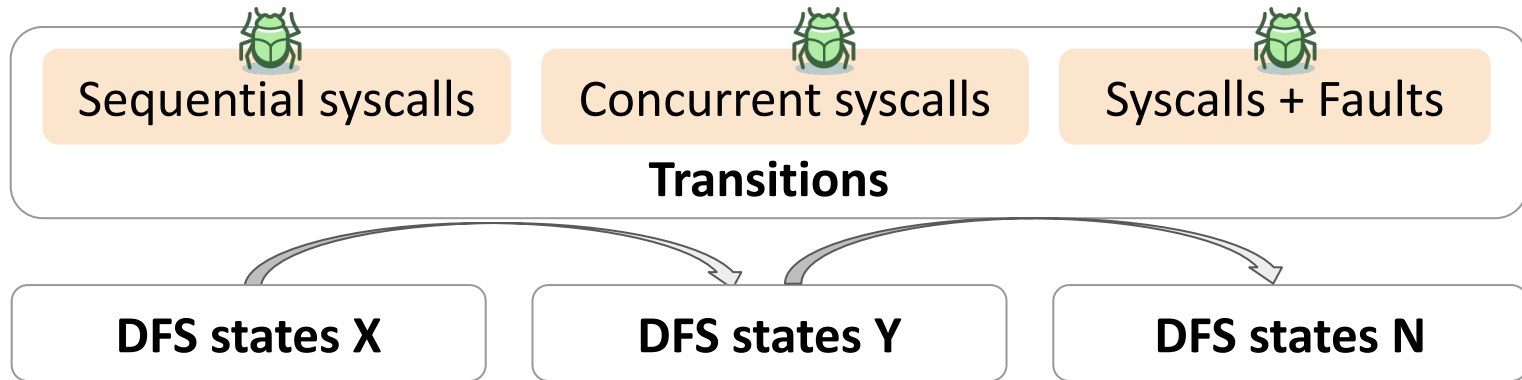
# What are the semantic bugs in DFS?

- **Semantic:** DFS state transitions specified in the spec
  - Syscalls: sequential and concurrent
  - Syscalls + faults



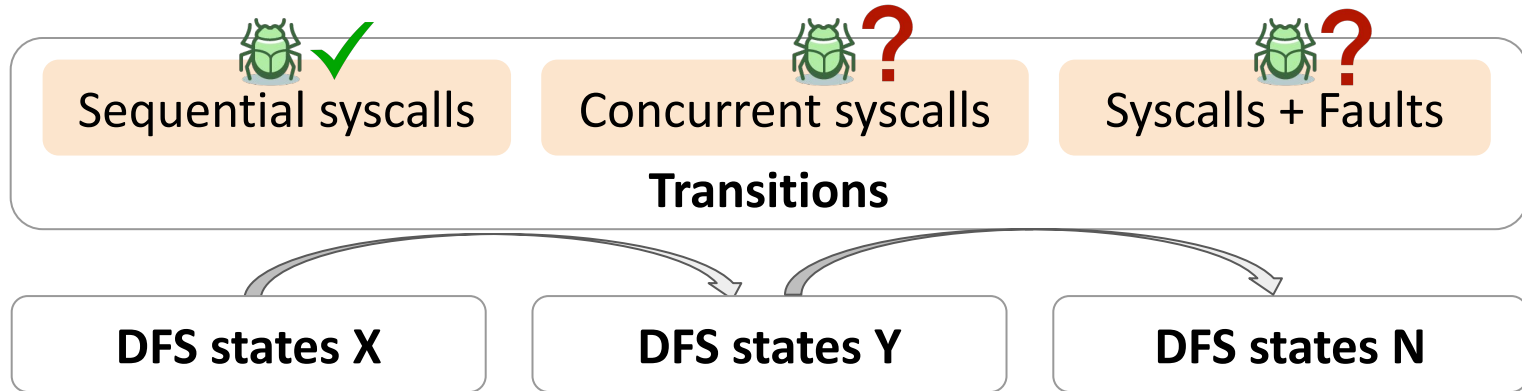
# What are the semantic bugs in DFS?

- **Semantic:** DFS state transitions specified in the spec
  - Syscalls: sequential and concurrent
  - Syscalls + faults
- **Semantic bugs:** implementations violate the specified semantic



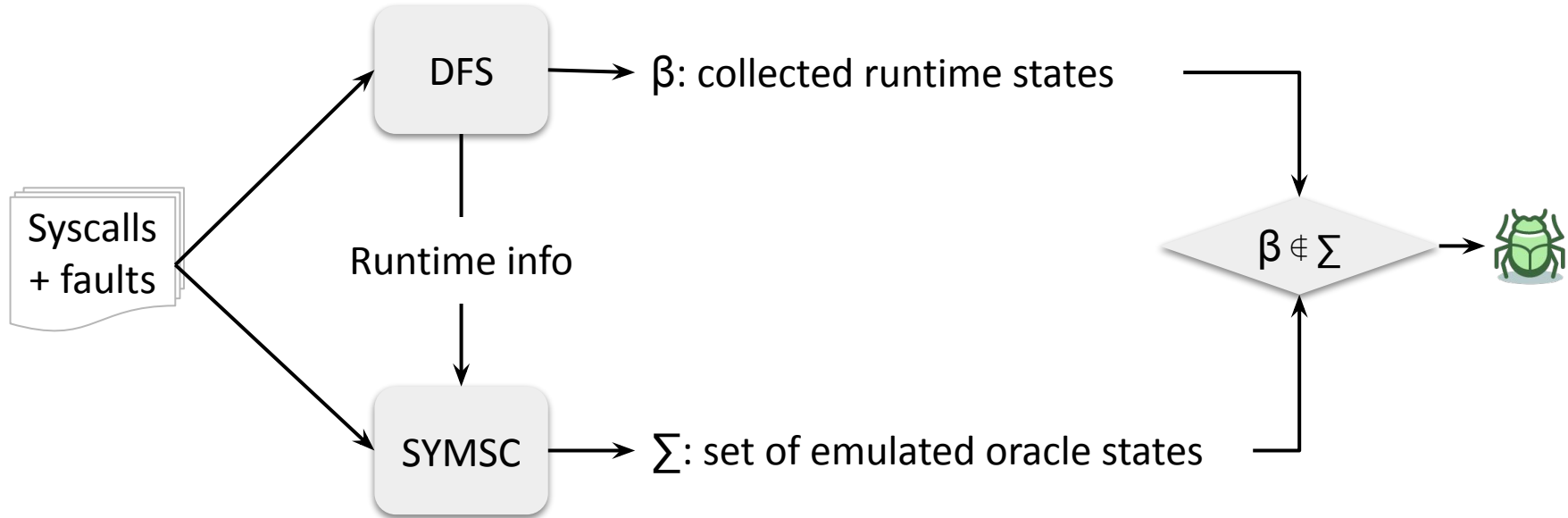
# Which semantic bug types are missing?

- **Semantic:** DFS state transitions specified in the spec
  - Syscalls: sequential and concurrent
  - Syscalls + faults
- **Semantic bugs:** implementations violate the specified semantic





# DFS semantic checker: SYMSC



# Emulate sequential syscall executions

- Emulate syscalls one by one according to the specification

## Sequential syscalls

mkdir A

setxattr A user.key:val

## Emulated states

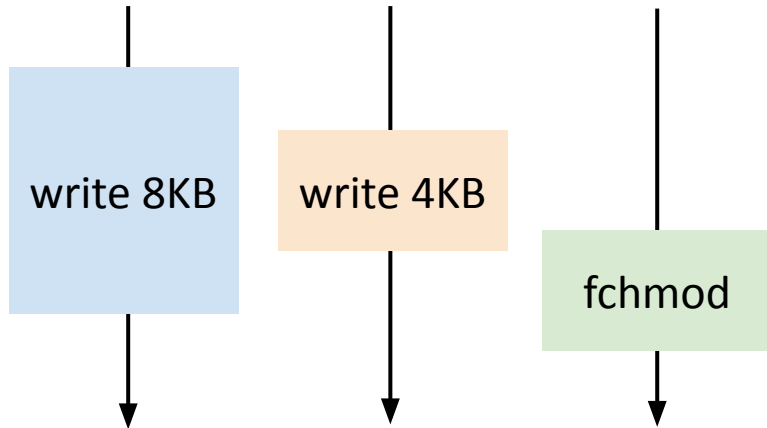
Inode0.dents = [., A]  
Inode1.dents = []

Inode1.xattr = [user.key:val]  
...

**But the spec does not specify semantics under concurrency and faults**

# How to emulate semantic under concurrency?

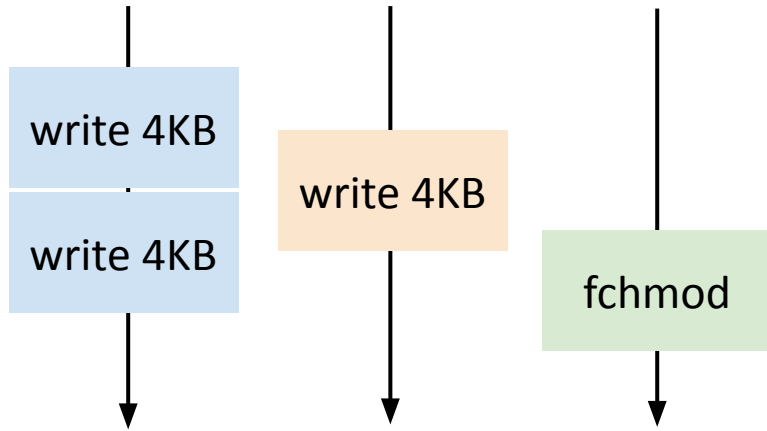
- Concurrency relation built from syscall's start and end timestamp



Concurrent execution

# Syscalls to atomic operations

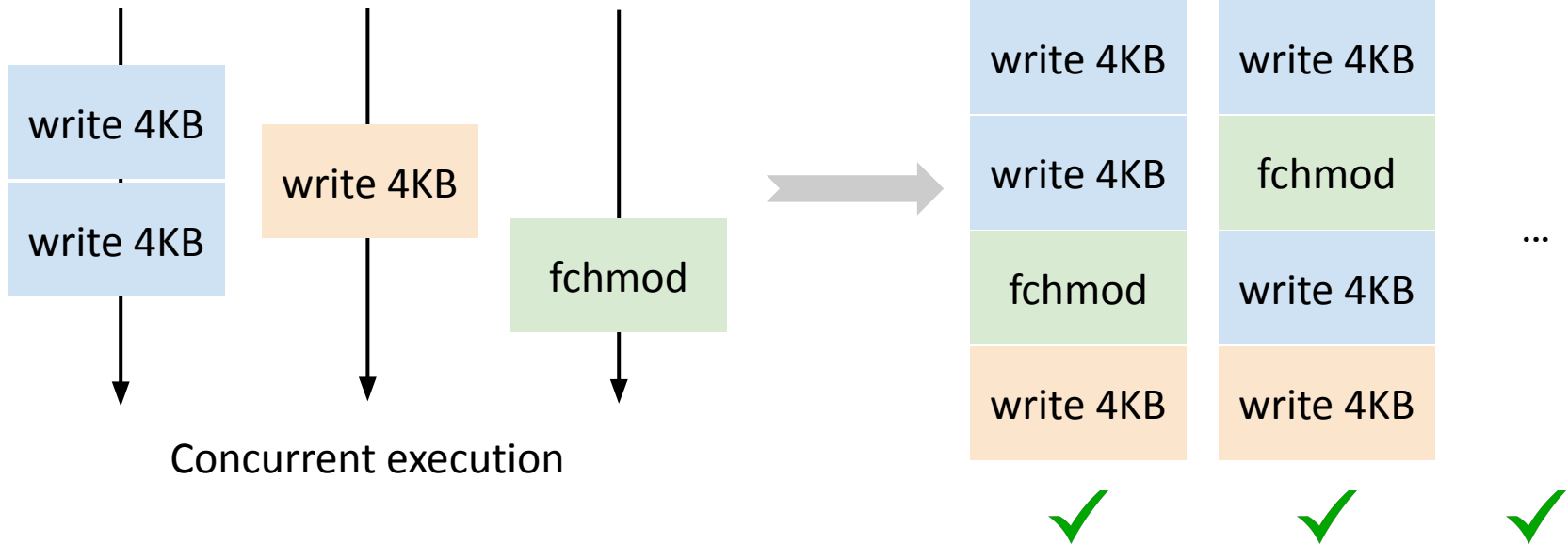
- Syscalls are splitted into one or more atomic operations



Concurrent execution

# Serialize concurrent atomic operations

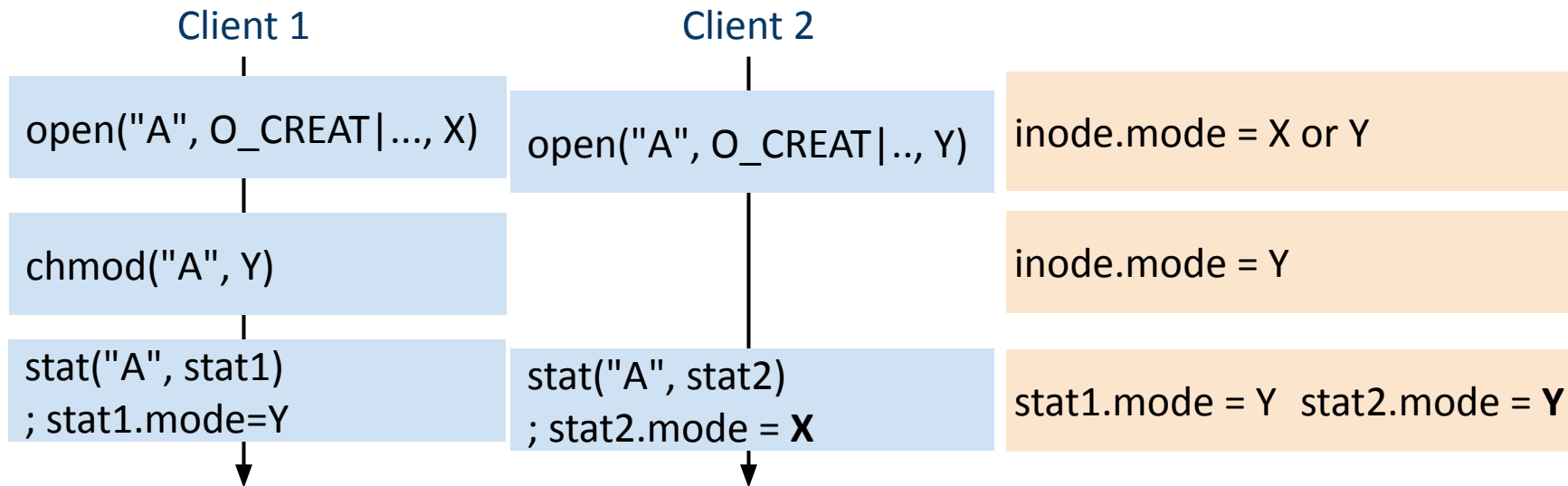
- Syscalls are splitted into one or more atomic operations
- **Oracle states:** emulate all interleavings among concurrent atomic operations



# Case study of a CephFS bug

## Runtime state

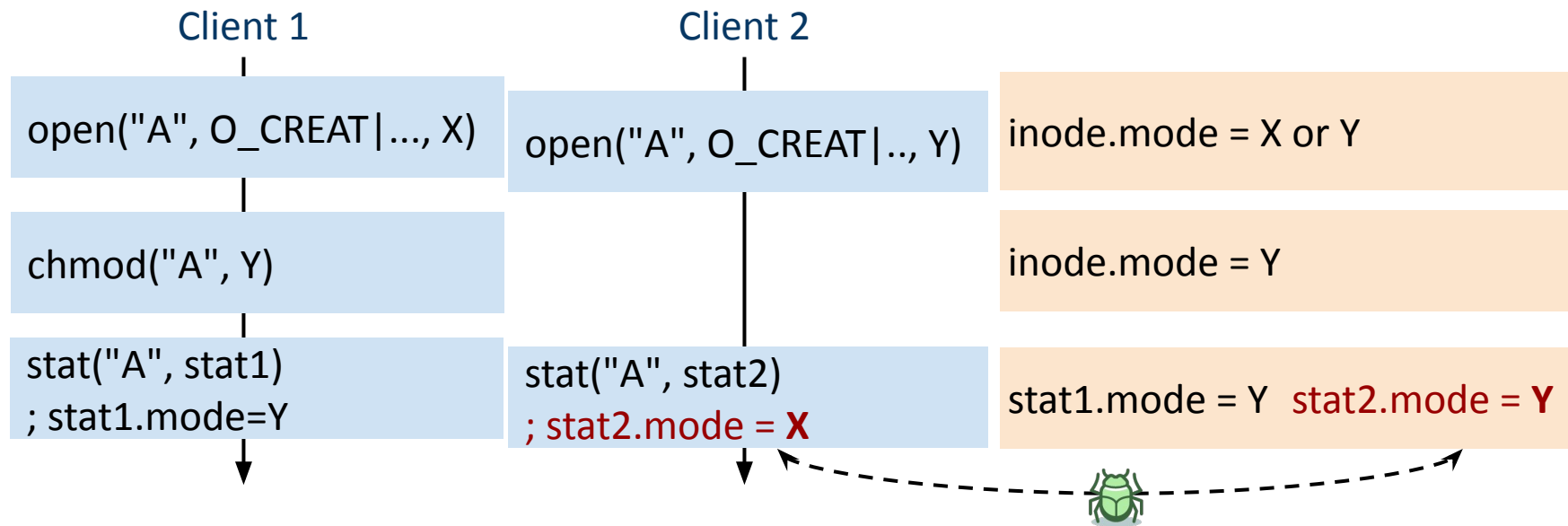
## Emulated states



# Case study of a CephFS bug

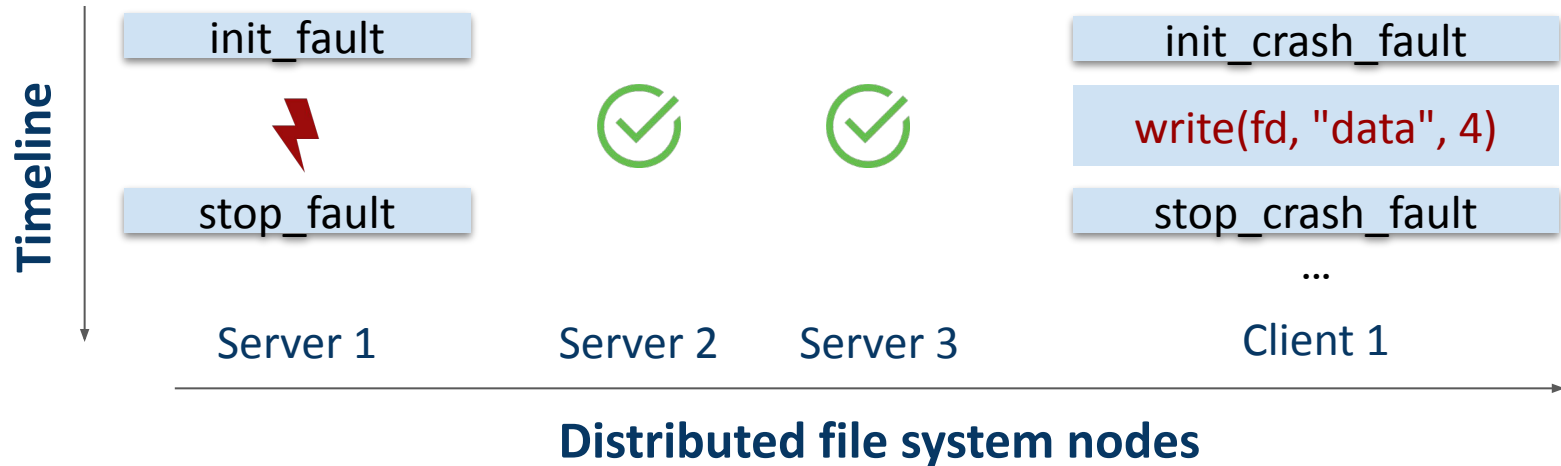
## Runtime state

## Emulated states



# What is the semantic under faults?

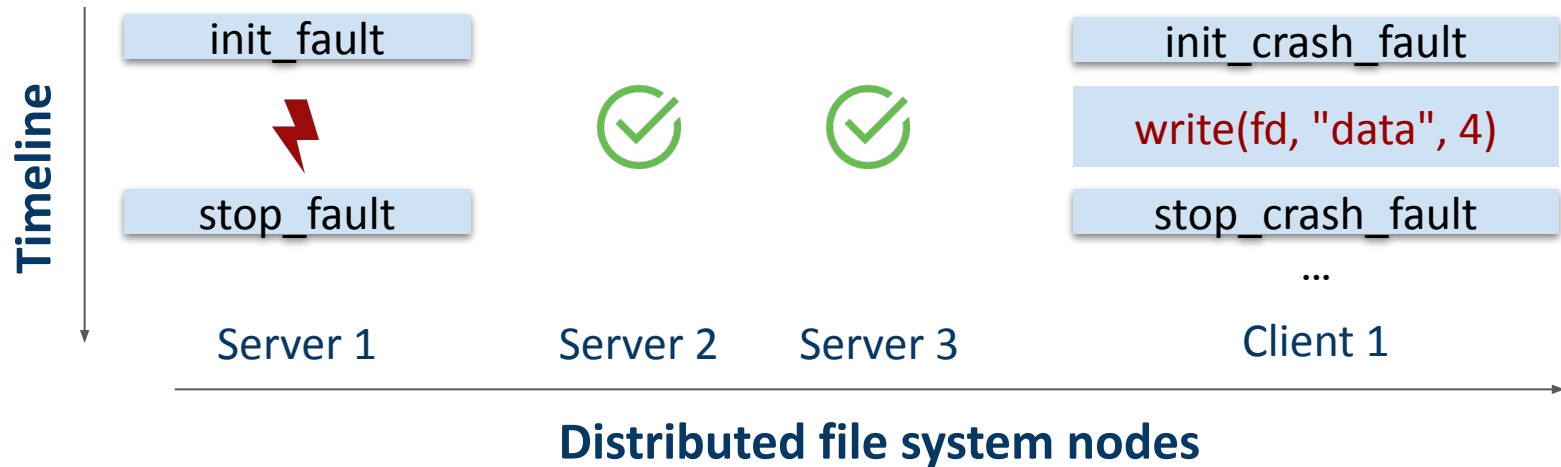
## What's the semantic of a write under faults?





# What is the semantic under faults?

- Apply the original semantic if the DFS is available to this syscall
- Return errors otherwise



# Implementation and evaluation setup

- Two 64-core machines
- Ran intermittently for two months on six targeted DFSes

lustre®









ORANGEFS



# How effective is Monarch in finding bugs?

- Monarch found bugs in all DFSes and categories, with a total of 48 bugs

DFS	Memory bugs	Semantic bugs
 Lustre	8	0
 GlusterFS	17	5
 OrangeFS	3	0
 BeeGFS	0	2
 CephFS	4	1
 NFS	8	0
<b>Total</b>	<b>40</b>	<b>8</b>

# What are the characteristics of DFS bugs?

**Faults play a critical role in exposing these bugs**

→ 14/40 memory bugs and 3/8 semantic bugs are exposed under faults

**Vulnerable code is scattered in both servers and clients**

→ 17 bugs (servers) vs 31 bugs (clients)

→ The root causes of semantics bugs are mostly in DFS servers

**Bug exposure might depend on specific DFS configurations**

# Conclusion



- **Problem:** Automatic and scalable bug-finding tool for DFSes
- **Our solution Monarch:** The first DFS fuzzing framework
  - Multi-node and cross-context fuzzing architecture
  - A two-step mutator to test DFSes with syscalls and faults
  - DFS semantic checker
- **Takeway:** Fuzzing is effective on distributed systems as well
- **Artifact:** <https://github.com/rs3lab/Monarch>

**Thank you!**