

Finding Crash-Consistency Bugs with Bounded Black-Box Testing

Jayashree Mohan, Ashlie Martinez, Soujanya Ponnappalli,
Pandian Raju, Vijay Chidambaram



vmware®



TEXAS

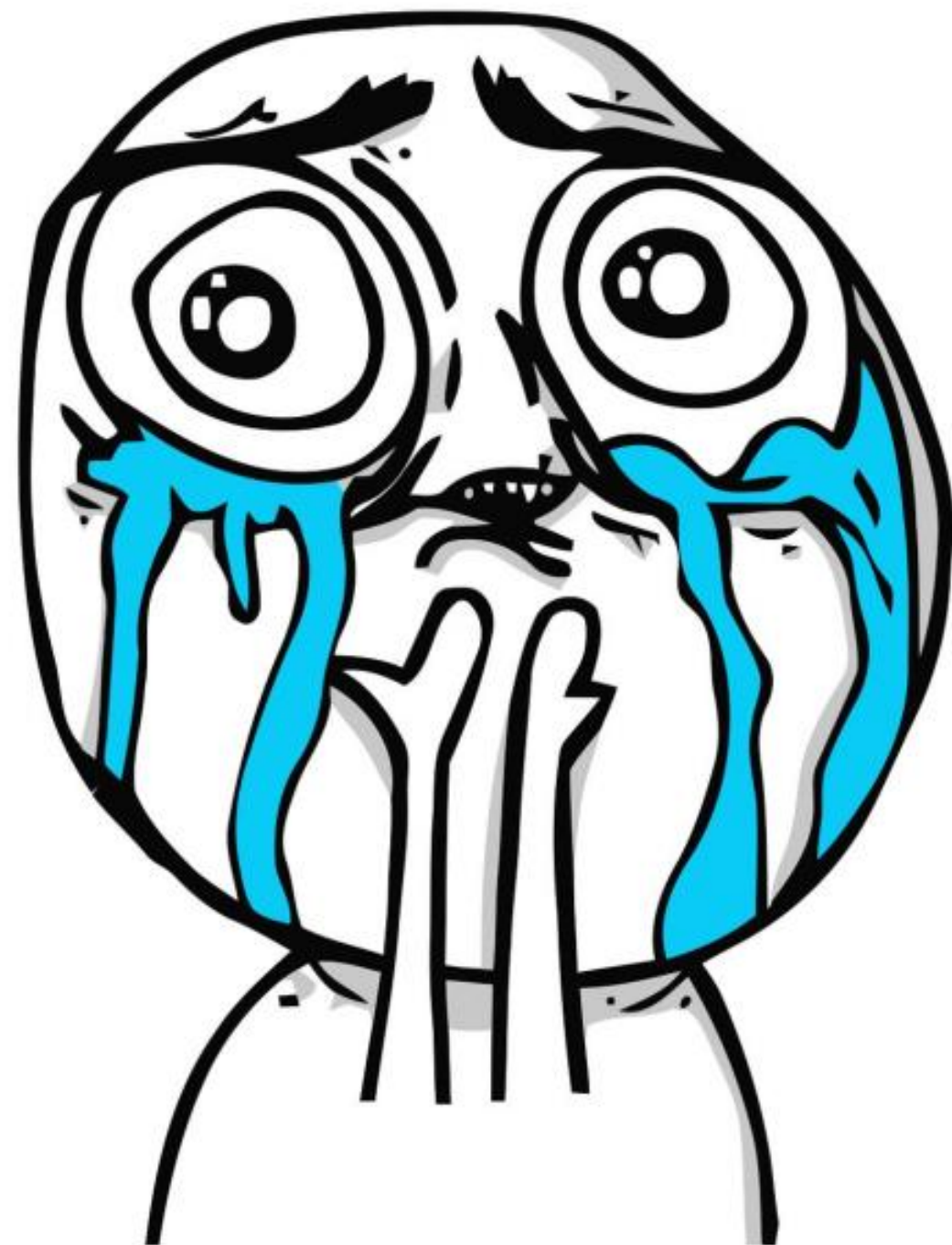
The University of Texas at Austin

Crashes



Image source : <https://www.fotolia.com>

**I wish filesystems
were crash-consistent!**

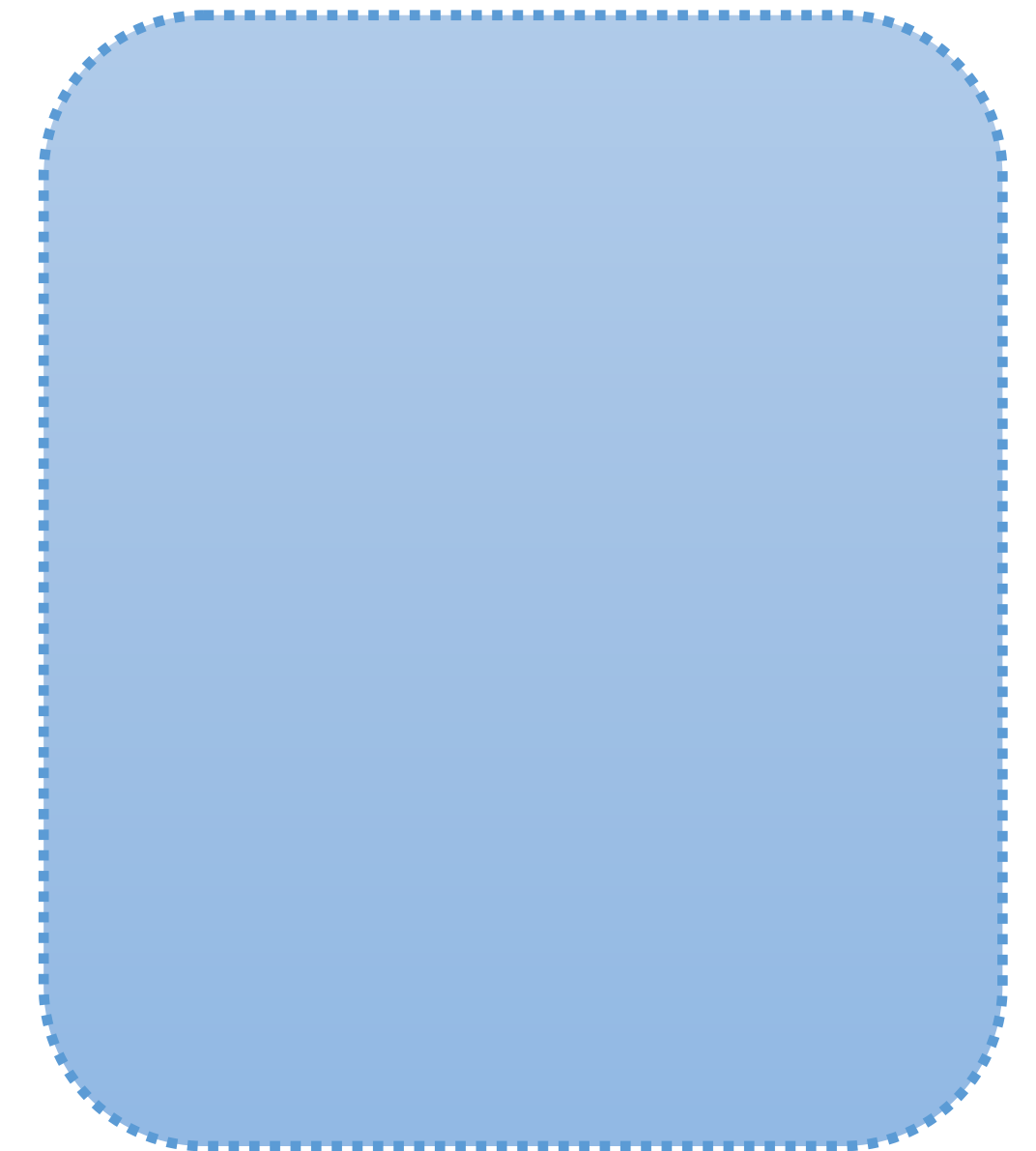


Rename atomicity bug in btrfs

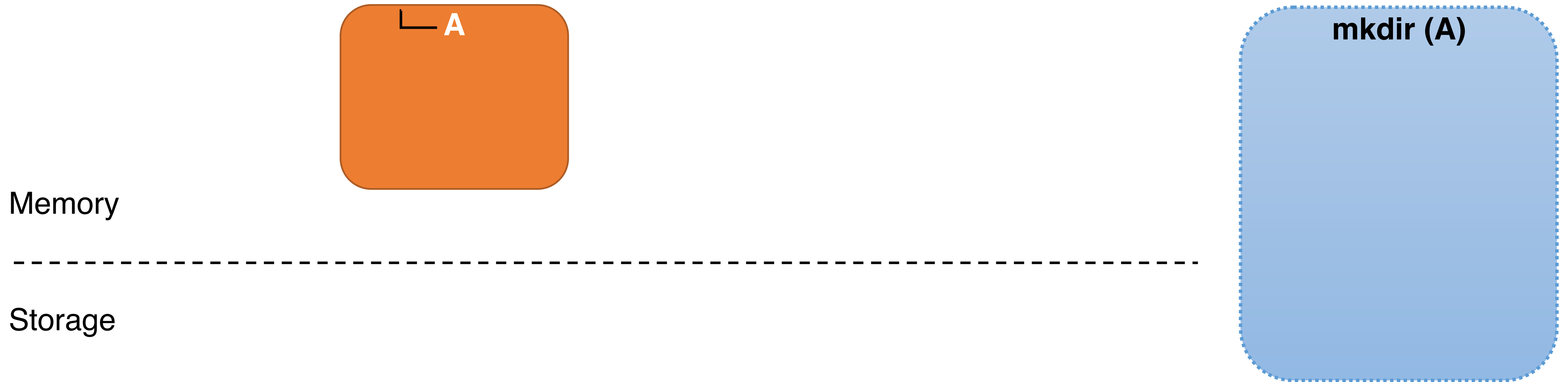
Memory



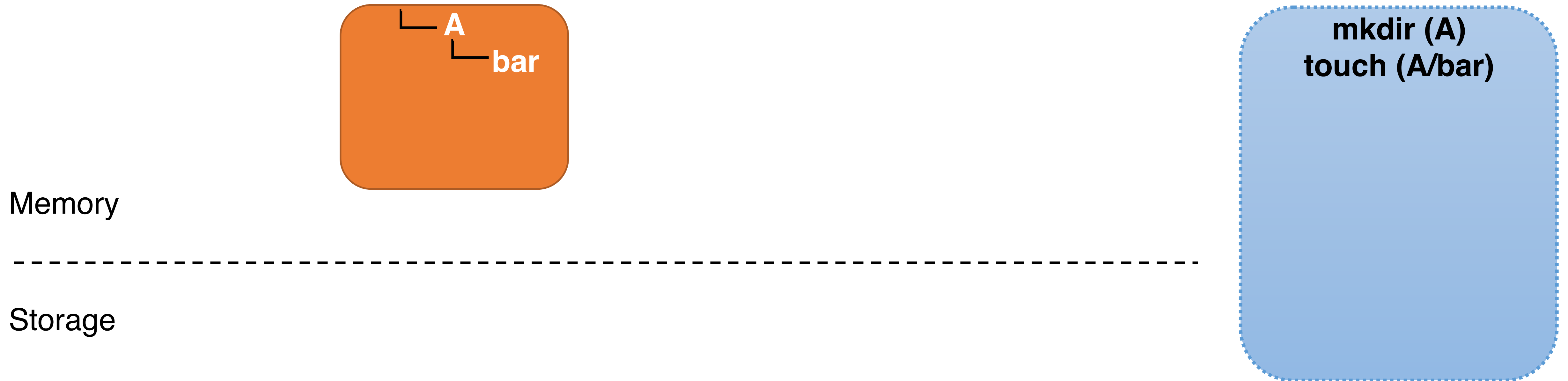
Storage



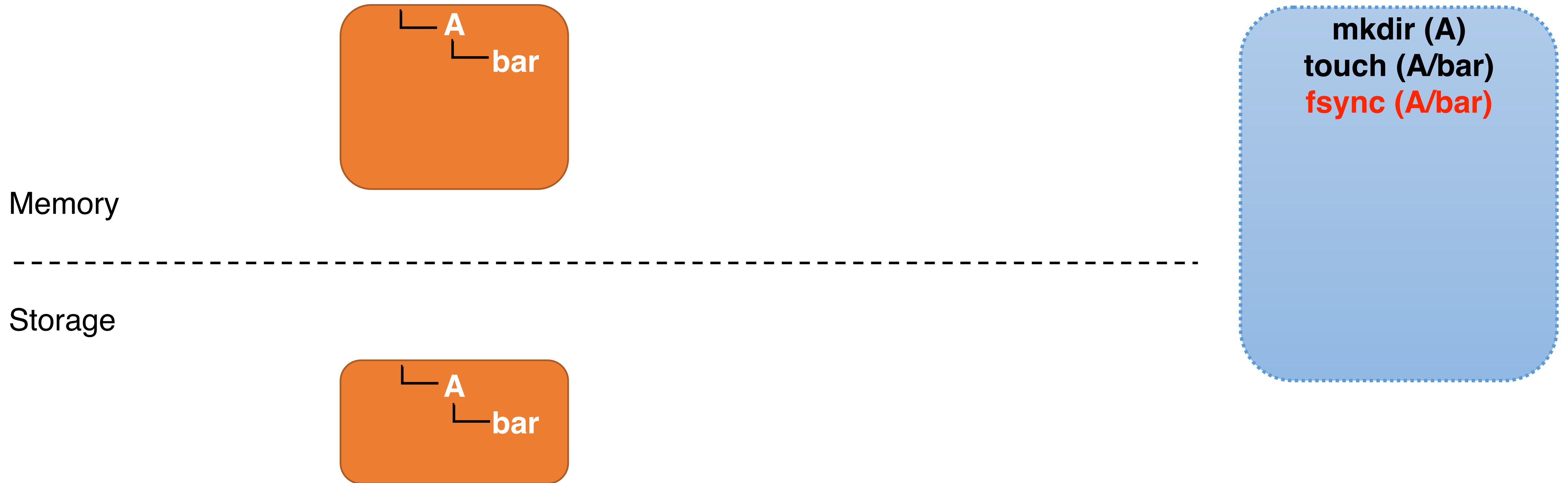
Rename atomicity bug in btrfs



Rename atomicity bug in btrfs

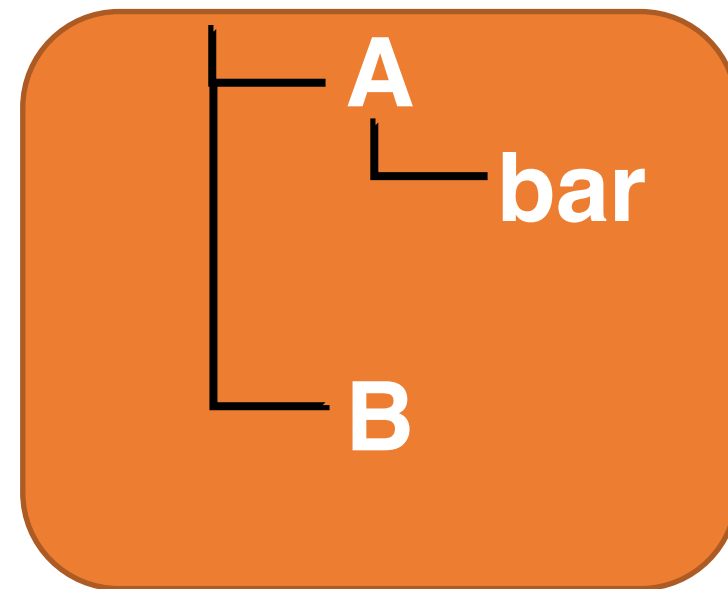


Rename atomicity bug in btrfs

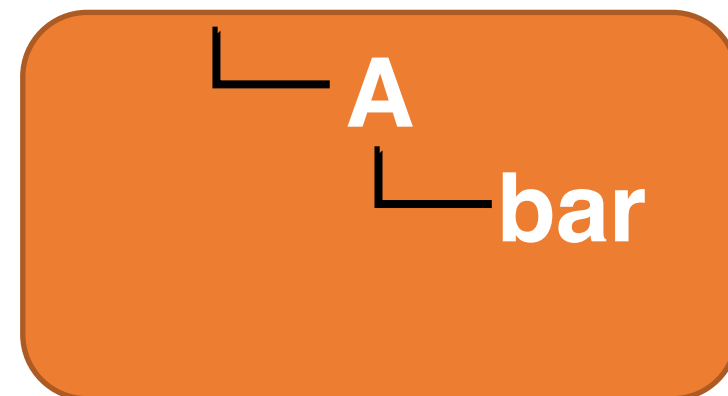


Rename atomicity bug in btrfs

Memory

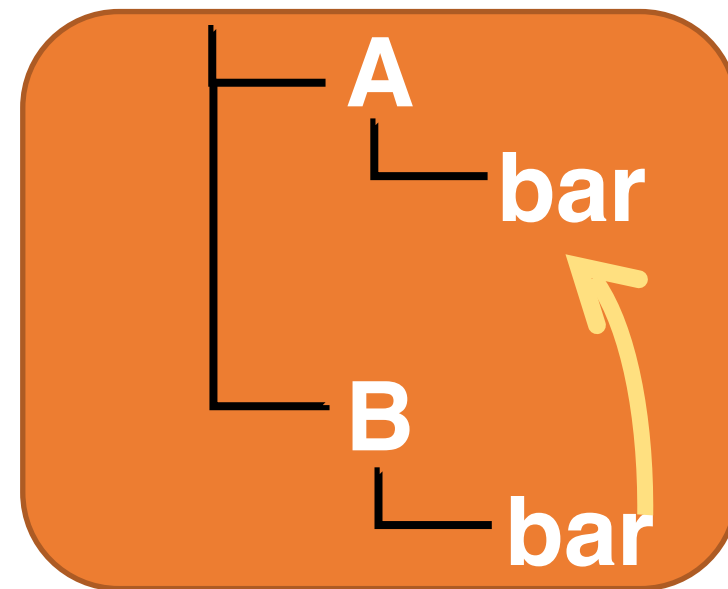


Storage

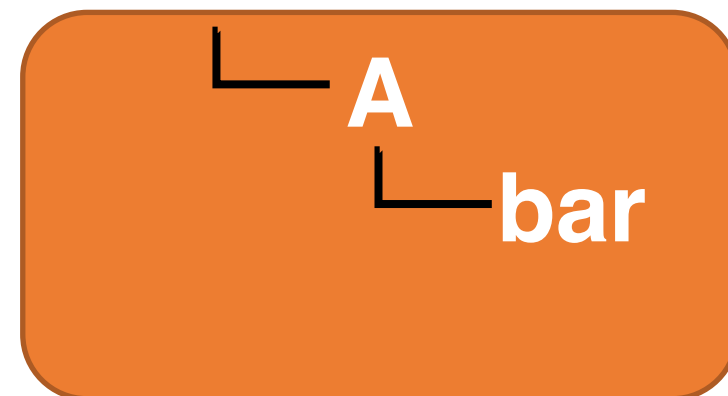


Rename atomicity bug in btrfs

Memory

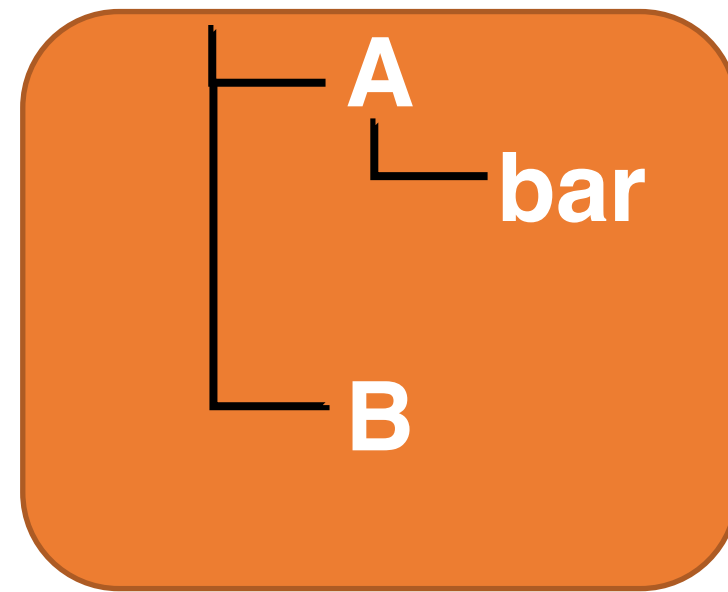


Storage

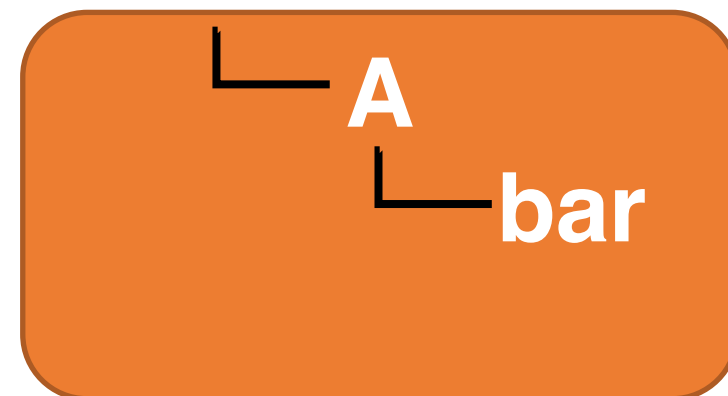


Rename atomicity bug in btrfs

Memory



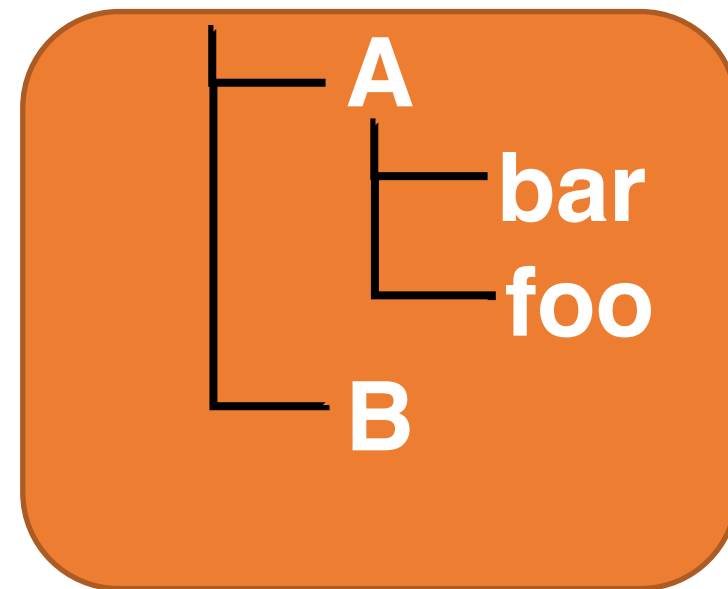
Storage



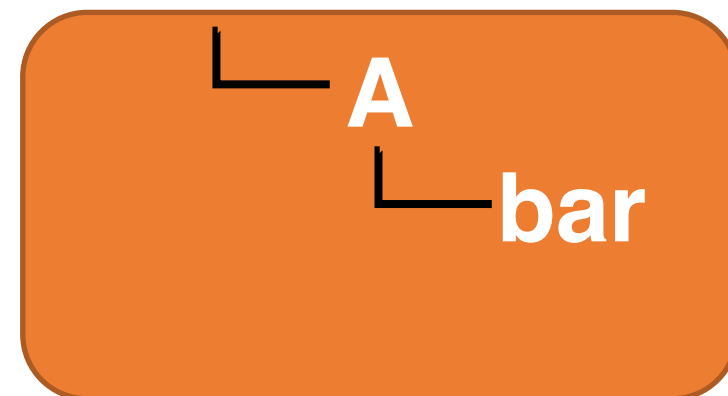
mkdir (A)
touch (A/bar)
fsync (A/bar)
mkdir (B)
touch (B/bar)
rename (B/bar, A/bar)

Rename atomicity bug in btrfs

Memory



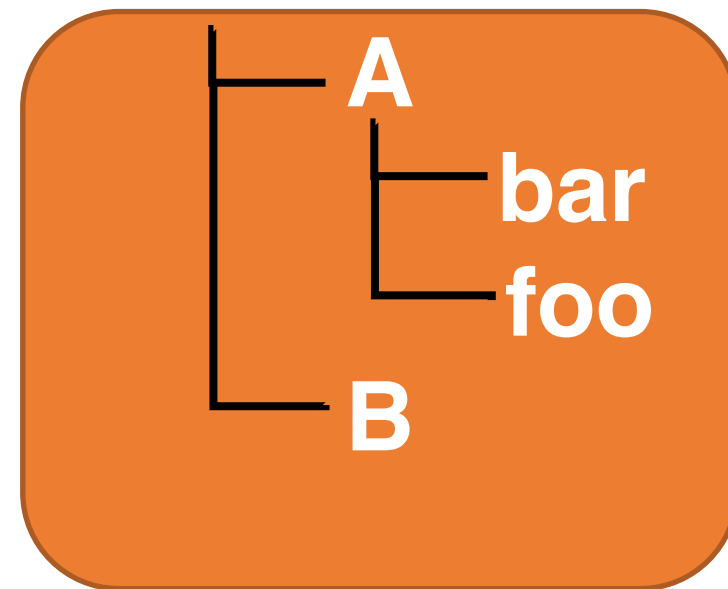
Storage



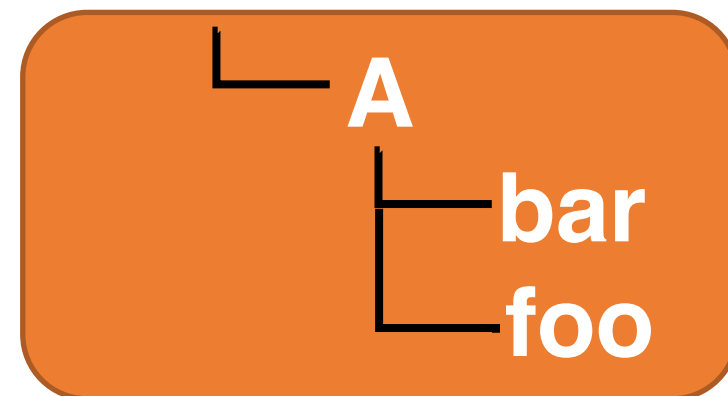
mkdir (A)
touch (A/bar)
fsync (A/bar)
mkdir (B)
touch (B/bar)
rename (B/bar, A/bar)
touch (A/foo)

Rename atomicity bug in btrfs

Memory



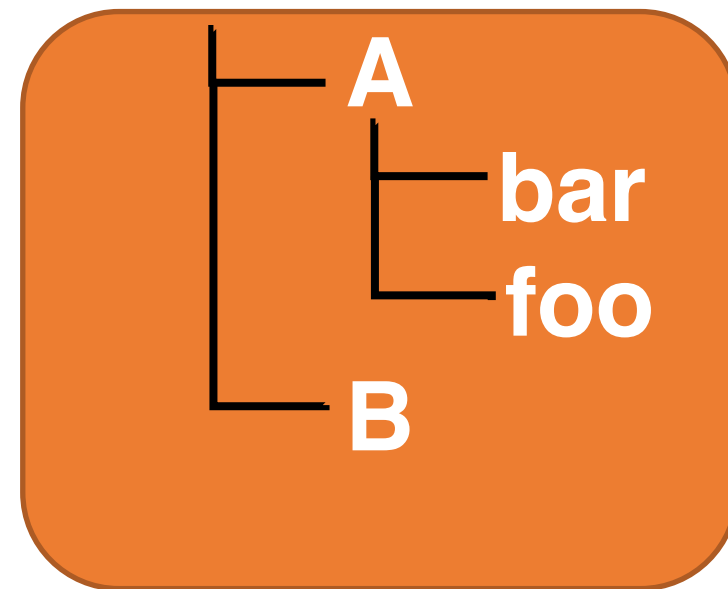
Storage



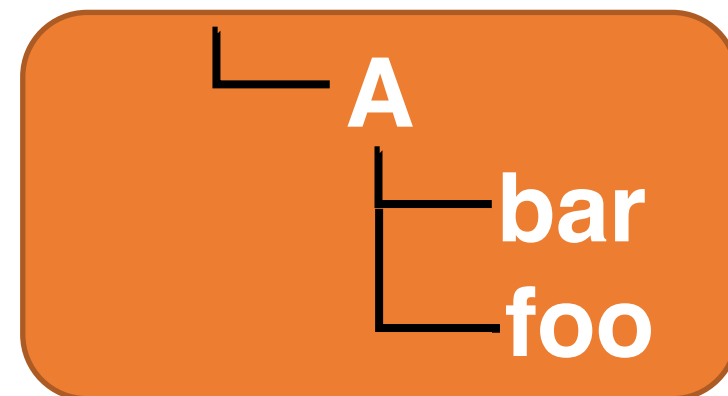
mkdir (A)
touch (A/bar)
fsync (A/bar)
mkdir (B)
touch (B/bar)
rename (B/bar, A/bar)
touch (A/foo)
fsync (A/foo)

Rename atomicity bug in btrfs

Memory



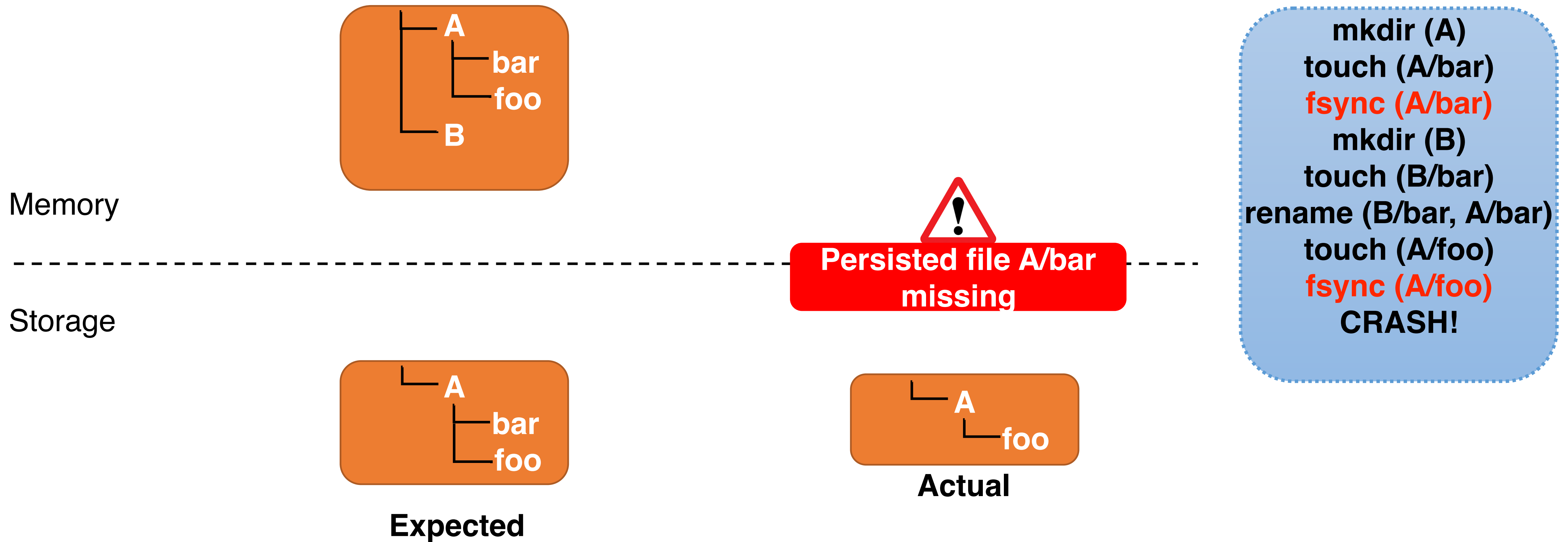
Storage



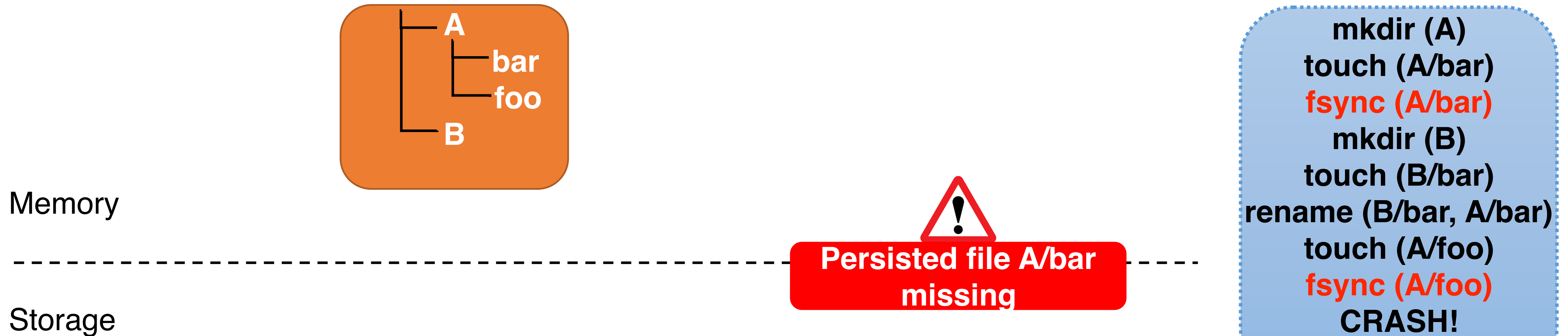
Expected

mkdir (A)
touch (A/bar)
fsync (A/bar)
mkdir (B)
touch (B/bar)
rename (B/bar, A/bar)
touch (A/foo)
fsync (A/foo)
CRASH!

Rename atomicity bug in btrfs



Rename atomicity bug in btrfs



**Exists in the kernel since 2014!
Found by ACE and CrashMonkey**

Testing Crash Consistency Today

- Build FS from scratch

**Verified
Filesystems**

Model Checking

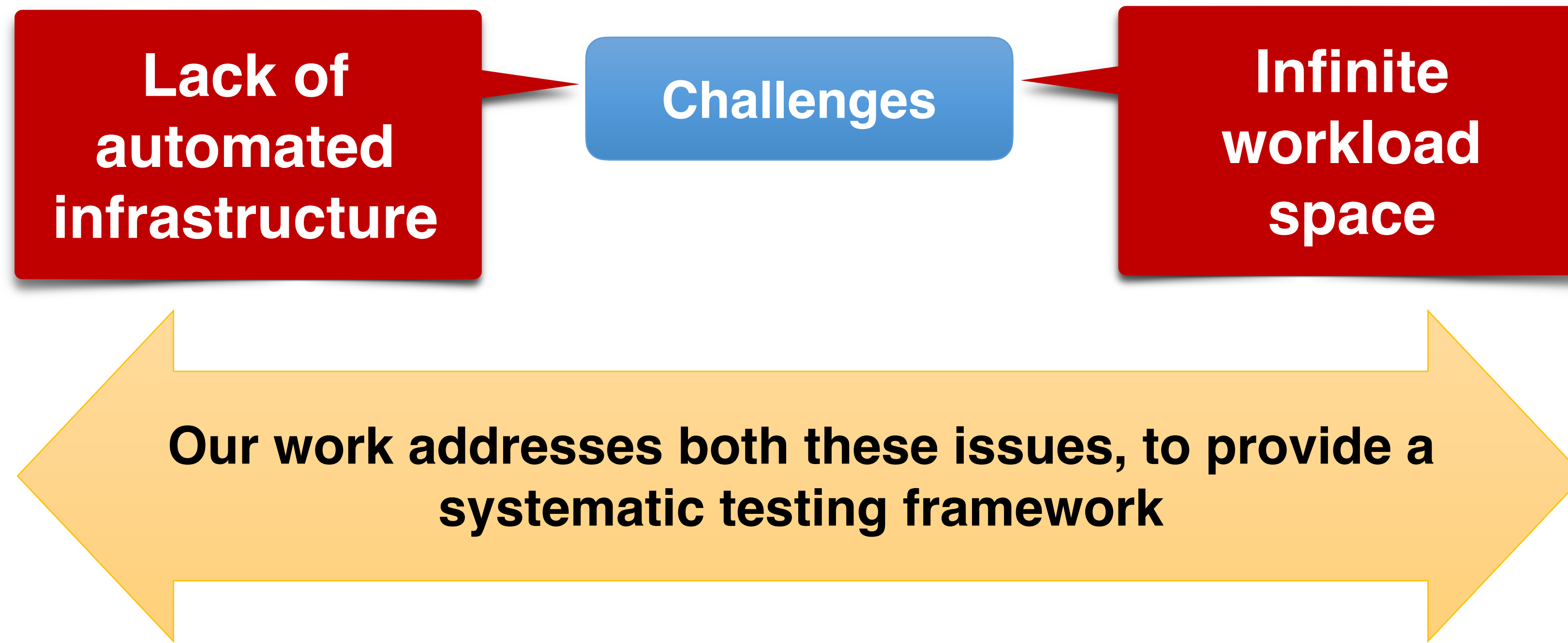
- Annotate filesystems
- Hard to do for existing FS

- State of the Art : xfstest suite
 - Collection of 482 regression tests

Only 5% of tests in xfstest check for file system crash consistency

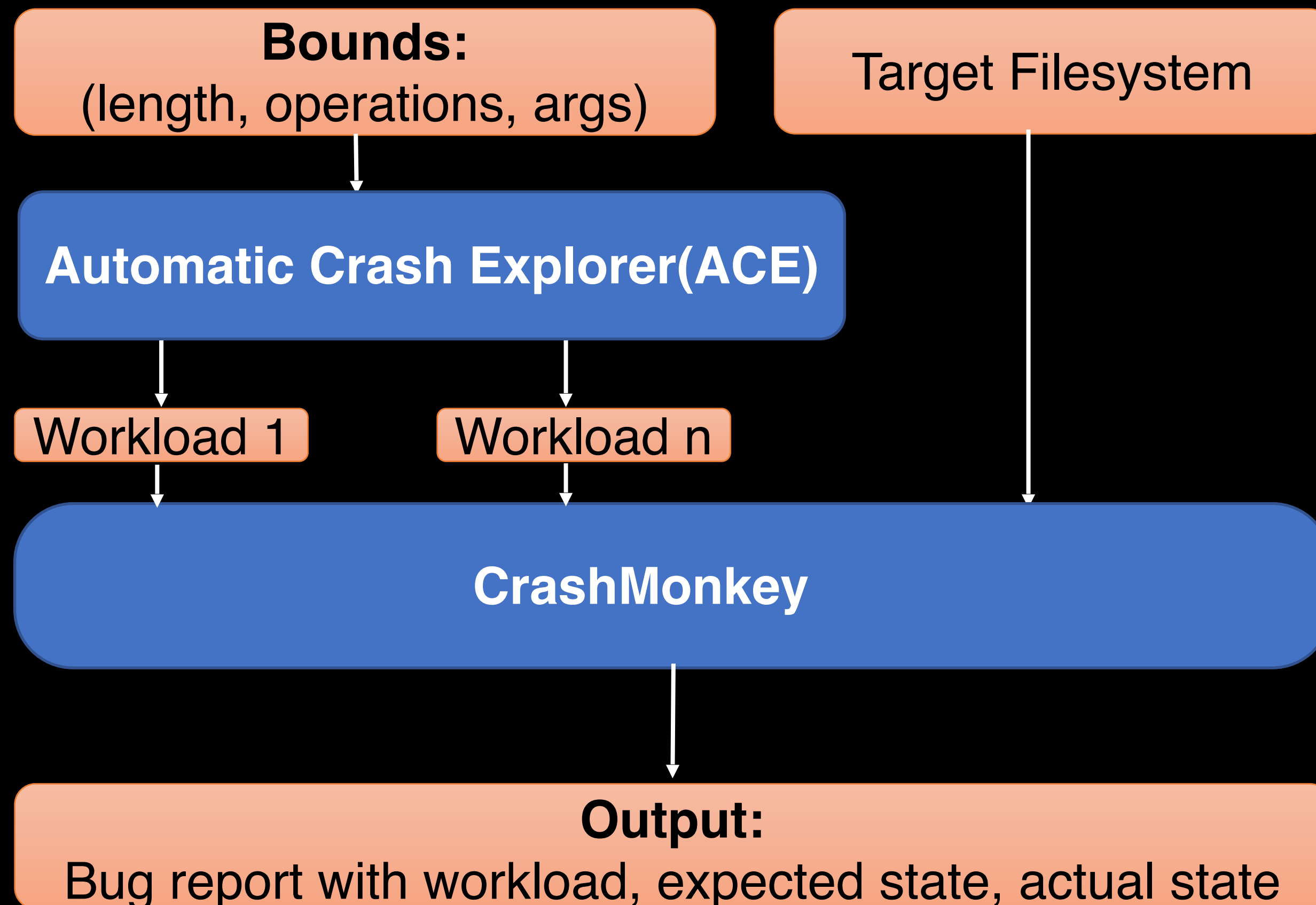
Challenges with systematic testing

Systematically generate workloads



Bounded Black-Box Crash Testing (B³)

New approach to testing file-system crash consistency

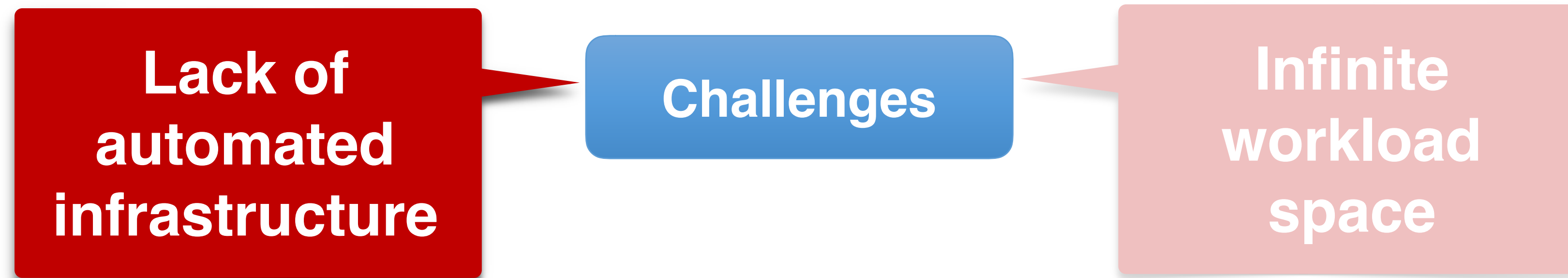


- ➔ Focus on reproducible bugs resulting in metadata corruption, data loss.
- ➔ Found 10 new bugs across btrfs and F2FS;
- ➔ Found 1 bug in FSCQ (verified file system)
- ➔ Filesystem agnostic – works with any POSIX file system

Outline

- CrashMonkey
- Bounded Black Box Crash Testing
- Automatic Crash Explorer (ACE)
- Demo

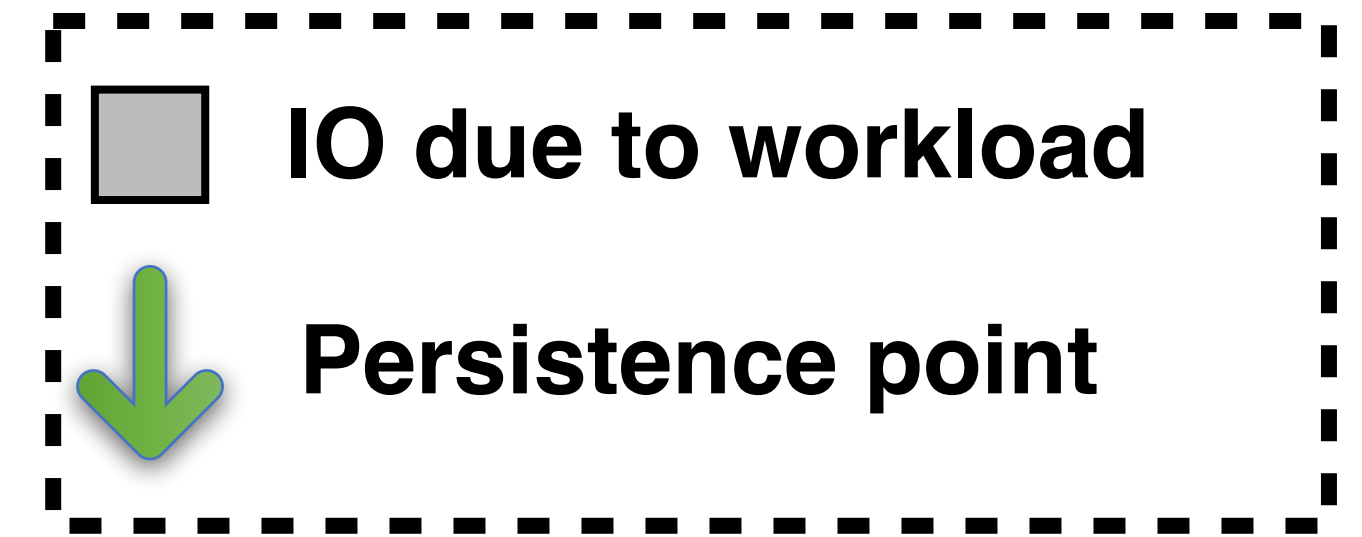
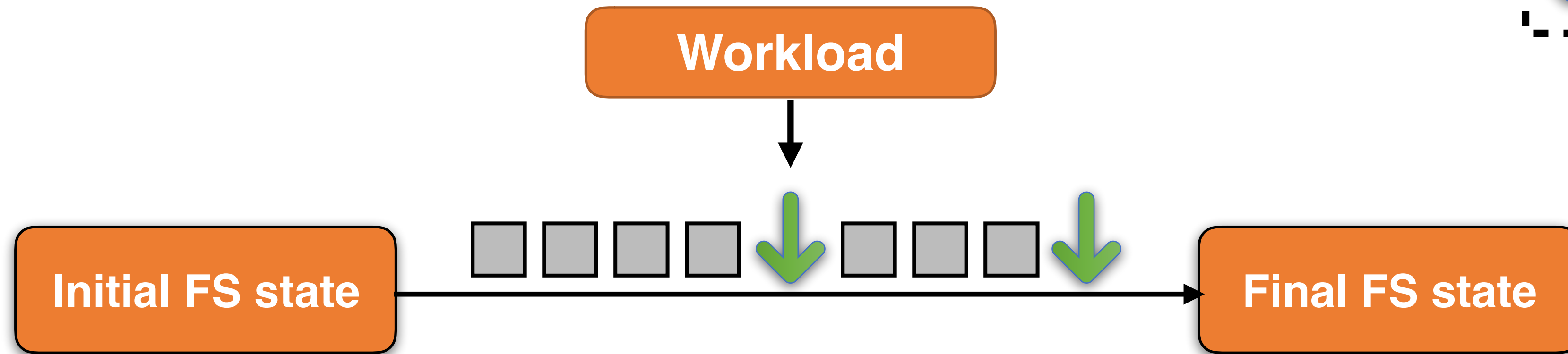
Challenges with systematic testing



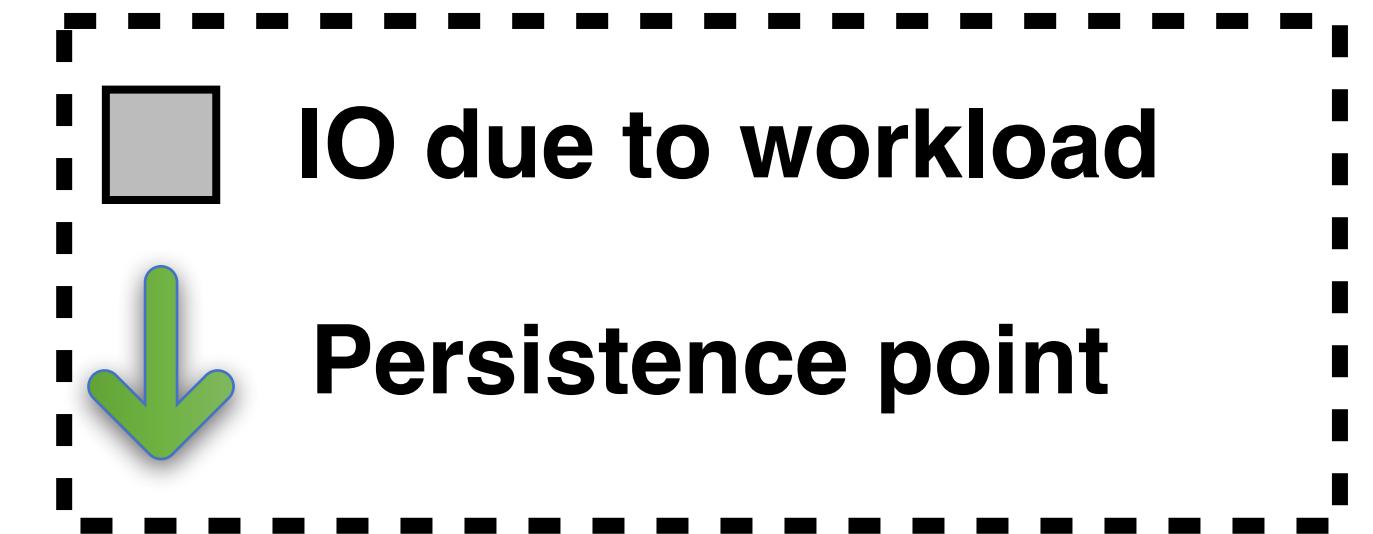
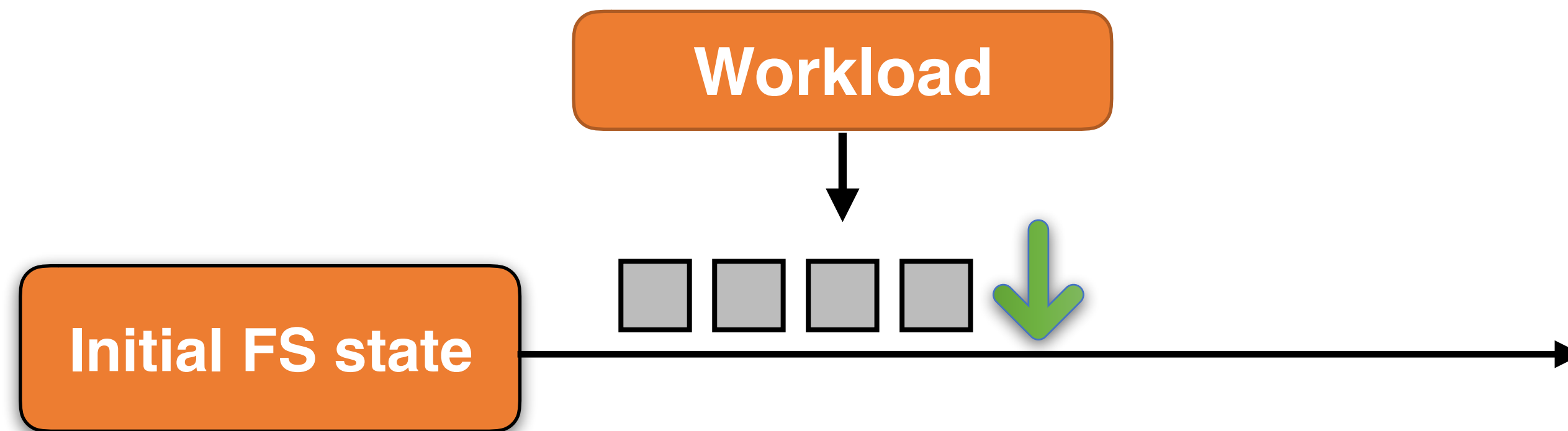
CrashMonkey

- Efficient infrastructure to record and replay block level IO requests
- Simulate crash at different points in the workload
- Automatically test for consistency after crash.
- Copy-on-write RAM block device

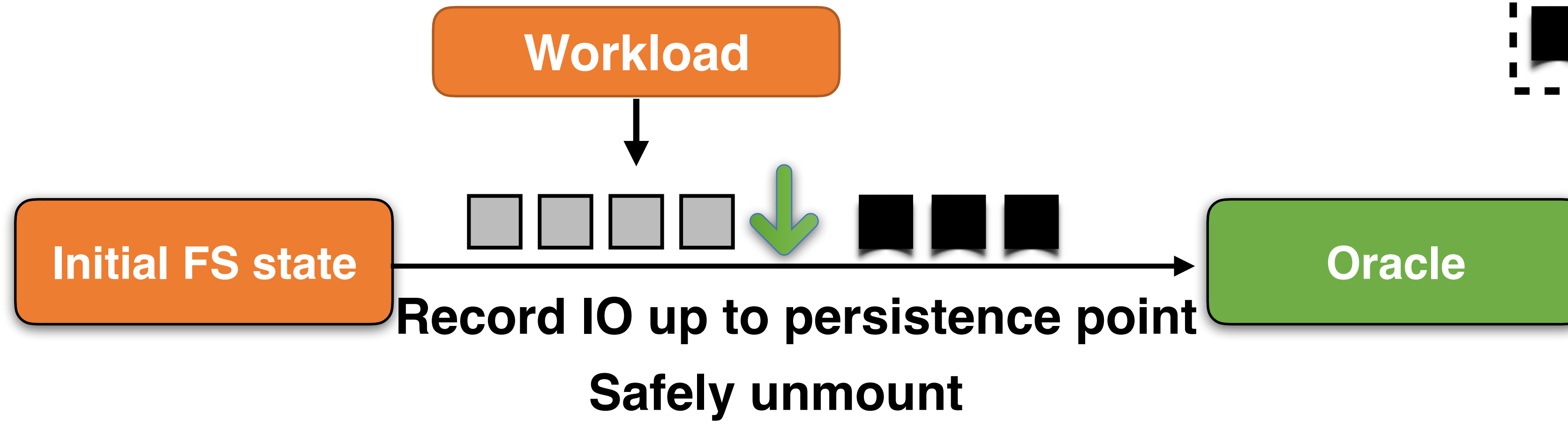
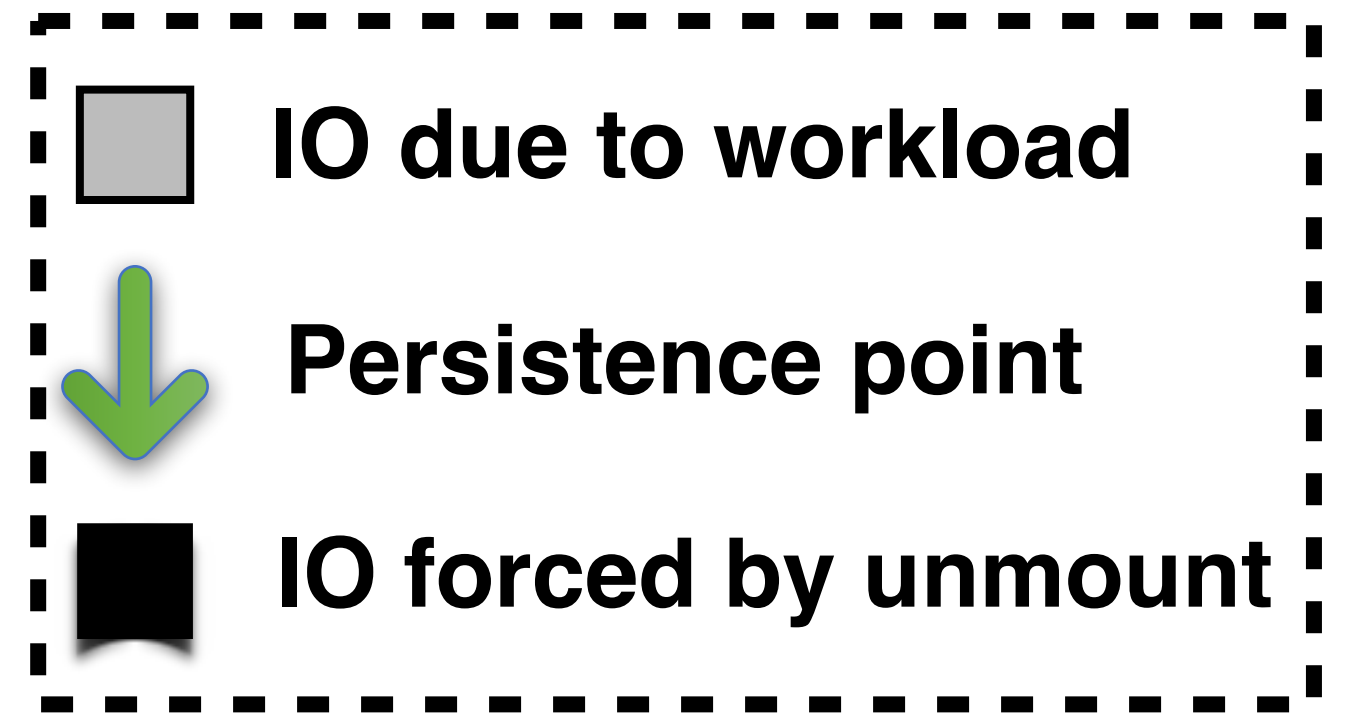
CrashMonkey in Action



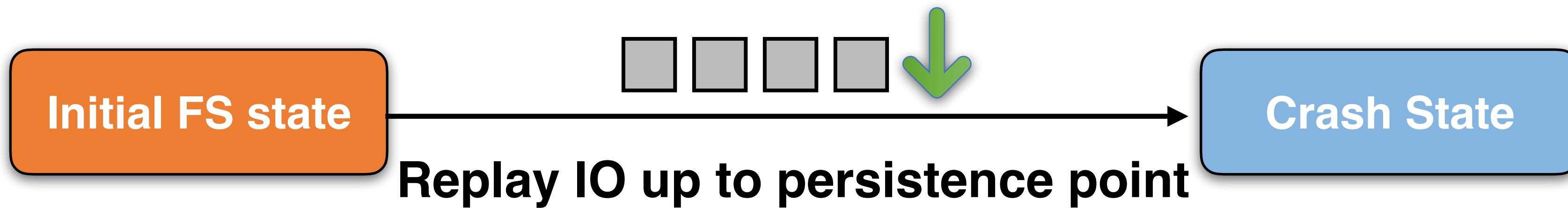
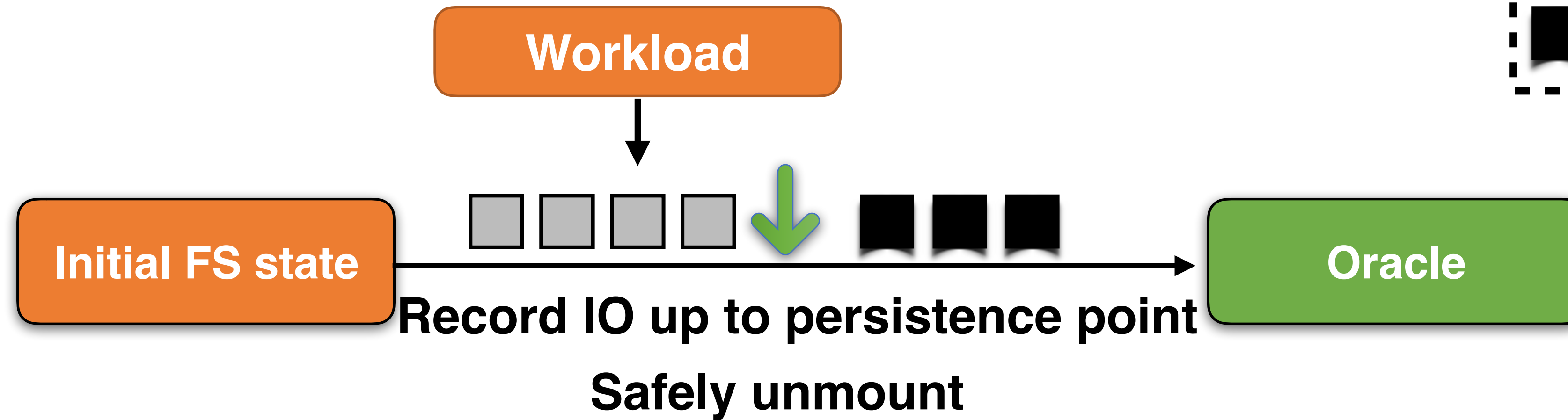
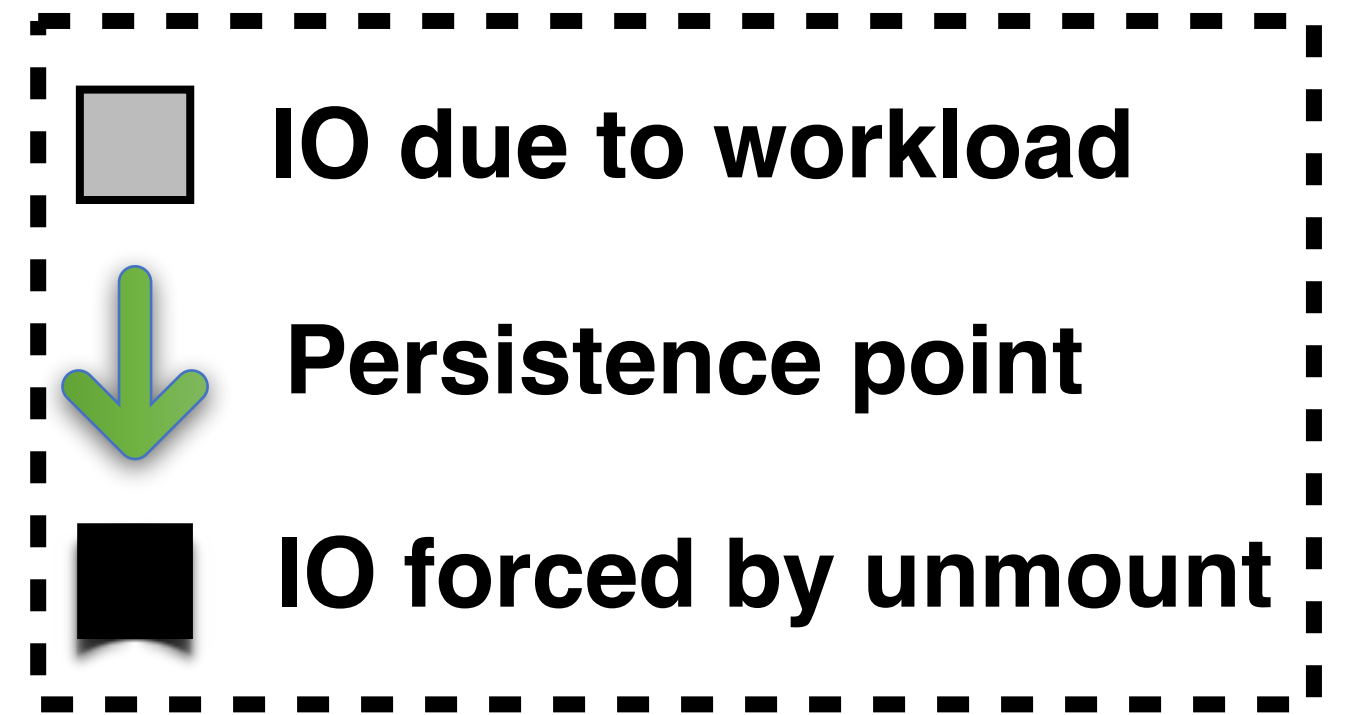
CrashMonkey in Action



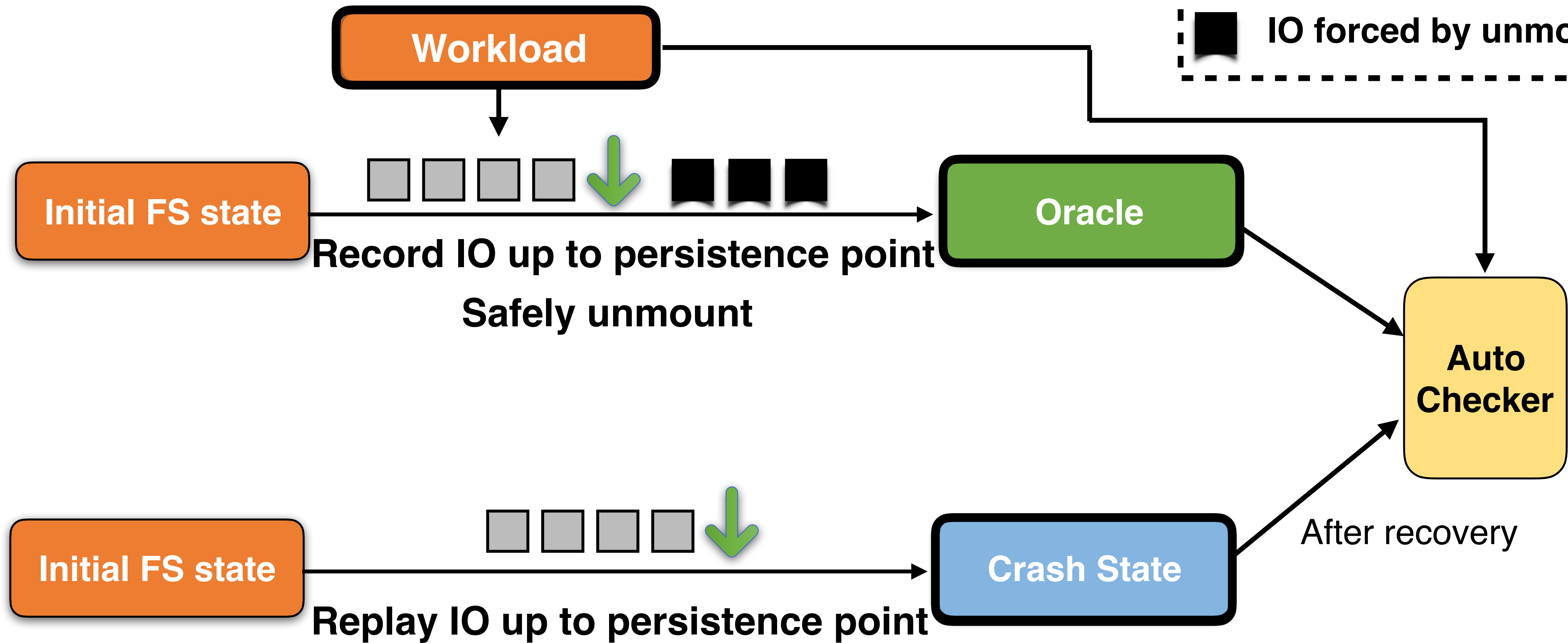
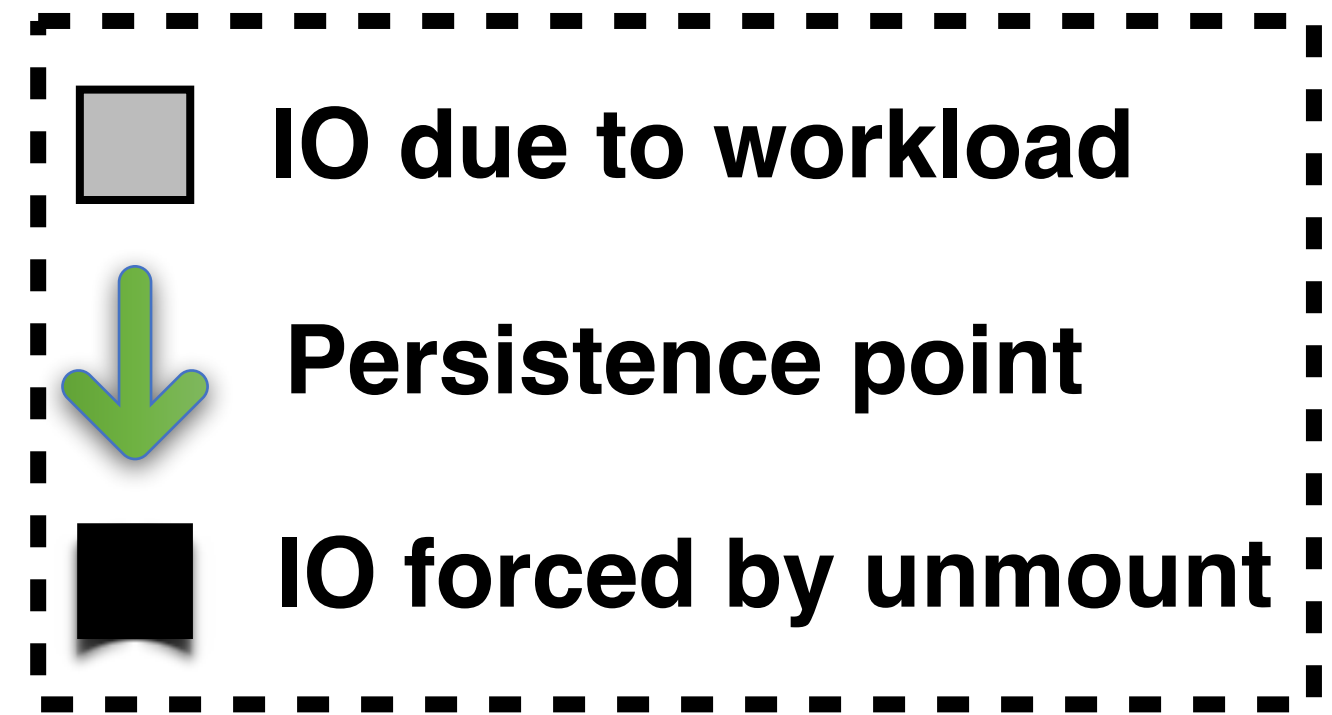
Phase 1 : Record IO



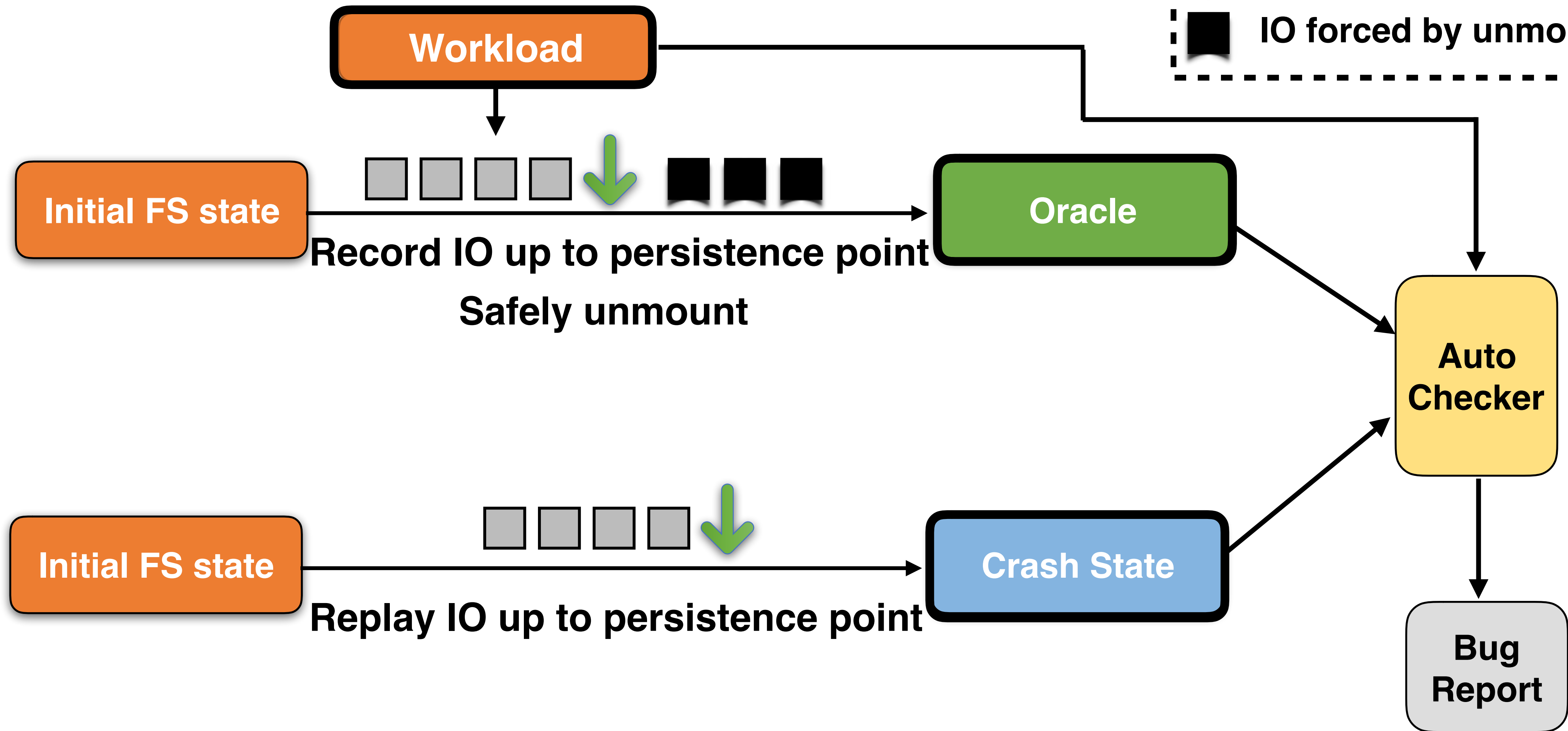
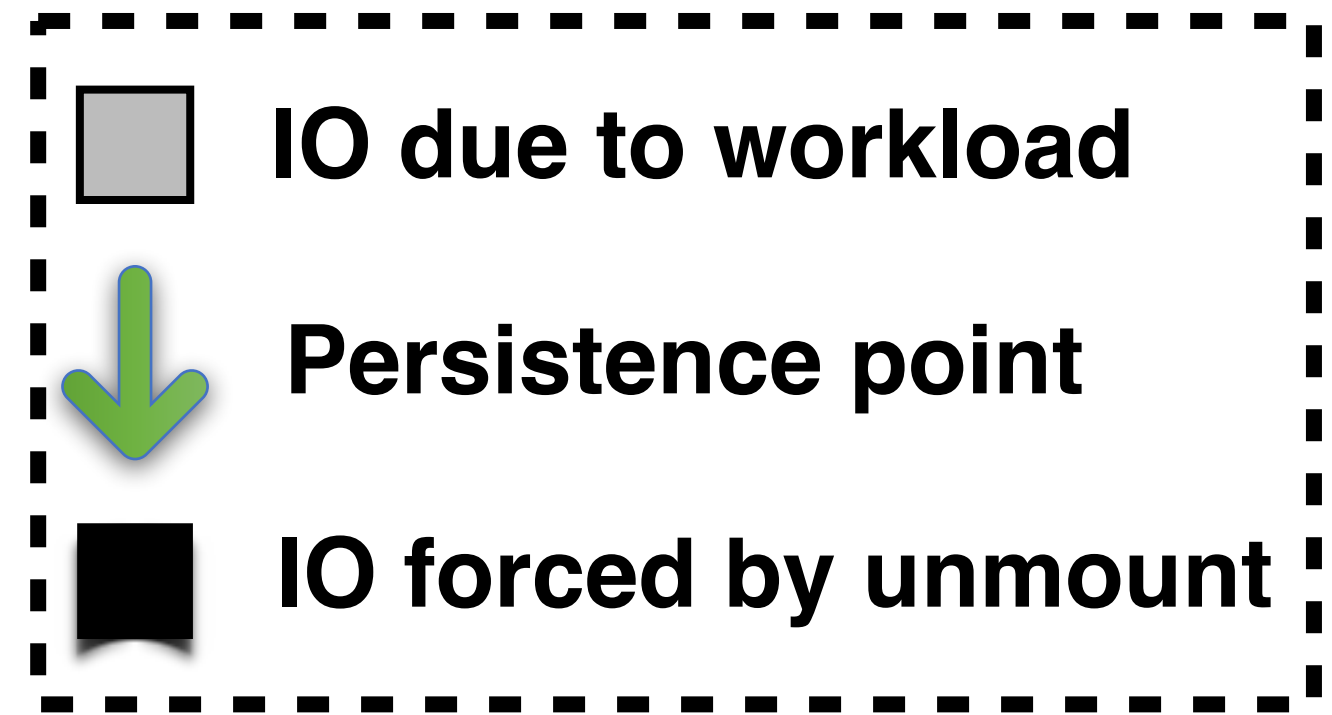
Phase 2 : Replay IO



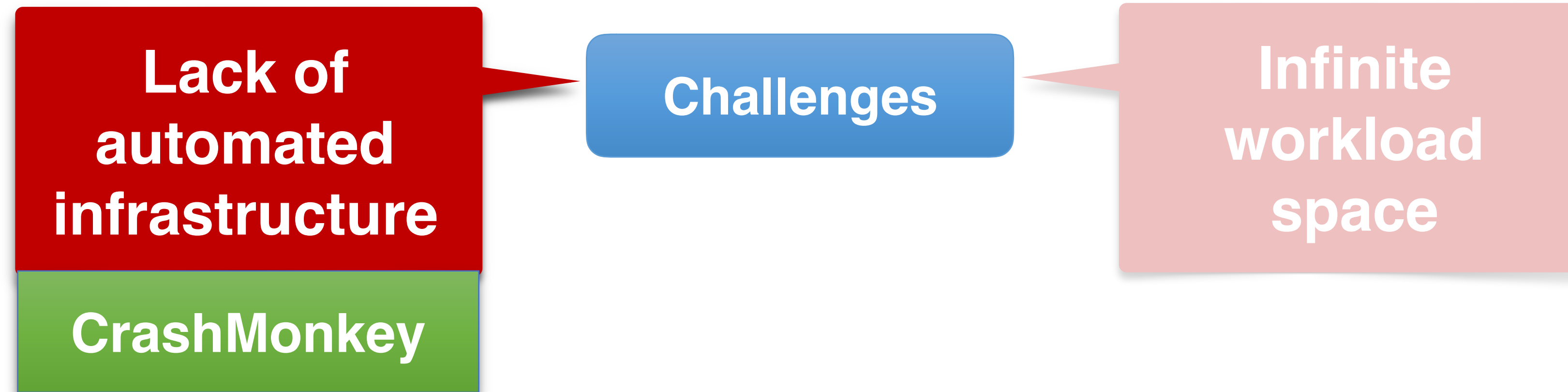
Phase 3 : Test for consistency



Phase 3 : Test for consistency



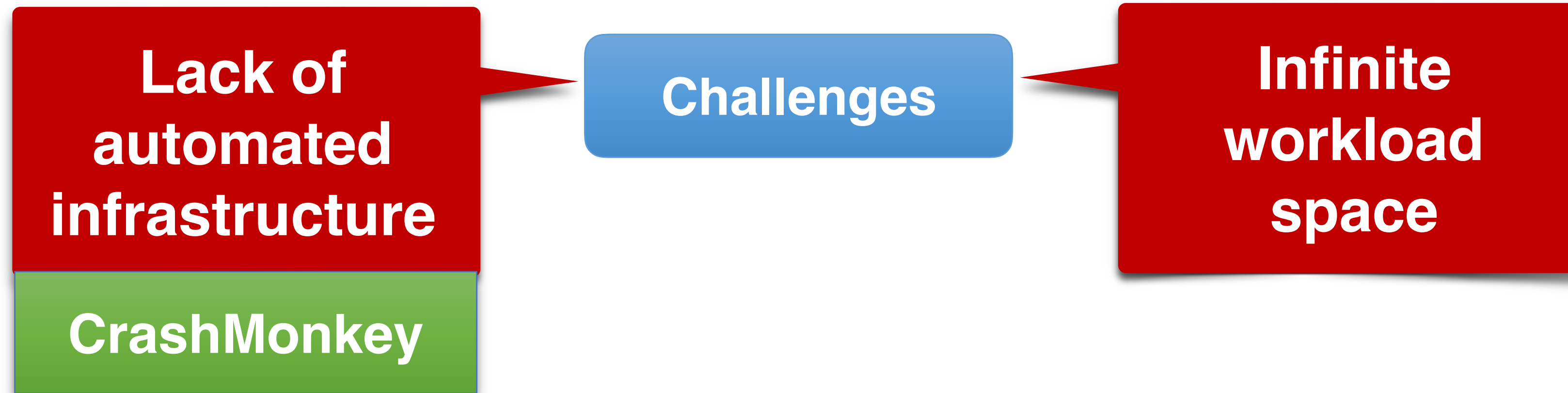
Challenges with Systematic Testing



So Far...

- Given a workload compliant to POSIX API, we saw how CrashMonkey generates crash states and automatically tests for consistency

Challenges with Systematic Testing



So Far...

- Given a workload compliant to POSIX API, we saw how CrashMonkey generates crash states and automatically tests for consistency
- Next question : How to automatically generate workloads in an the infinite workload space?

Exploring the infinite workload space

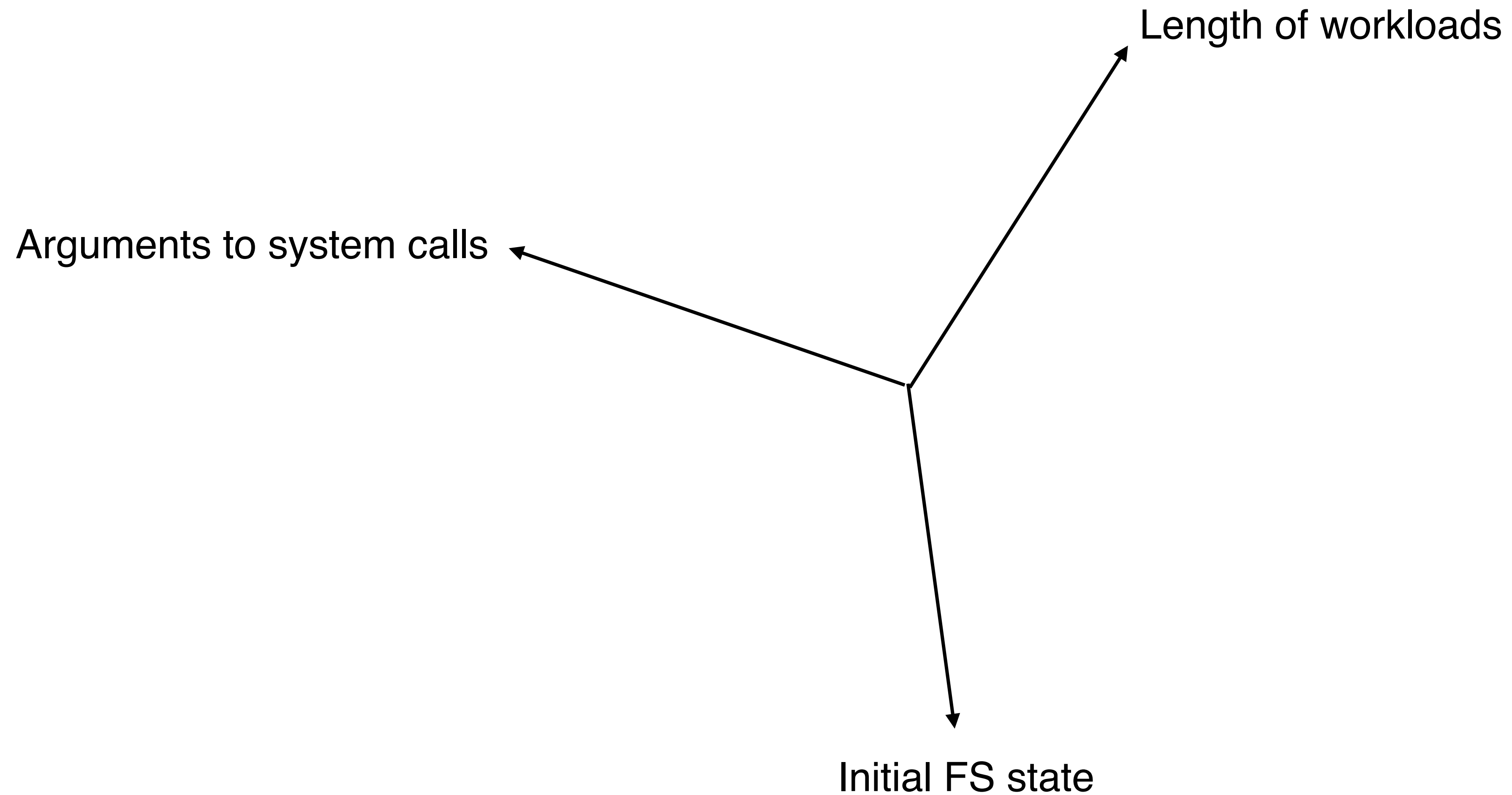
Challenges:

- Infinite length of workloads
- Large set of filesystem operations
- Infinite parameter options (file/directory names, depth)
- Infinite options for initial filesystem state
- When in the workload to simulate a crash?

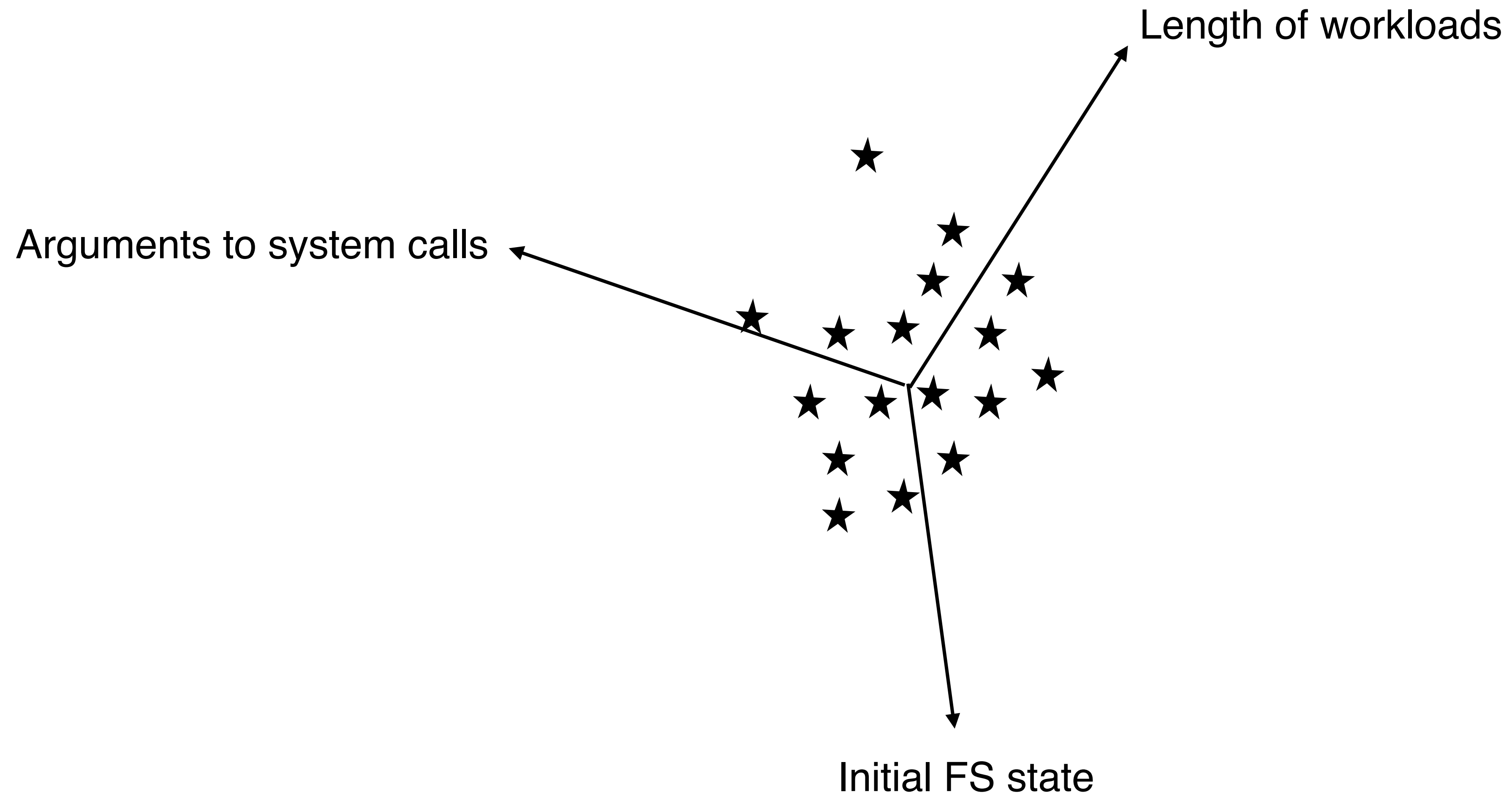
Outline

- CrashMonkey
- **Bounded Black Box Crash Testing**
- Automatic Crash Explorer (ACE)
- Demo

B³ : Bounded Black Box Crash Testing



B³ : Bounded Black Box Crash Testing



B³ : Bounded Black Box Crash Testing

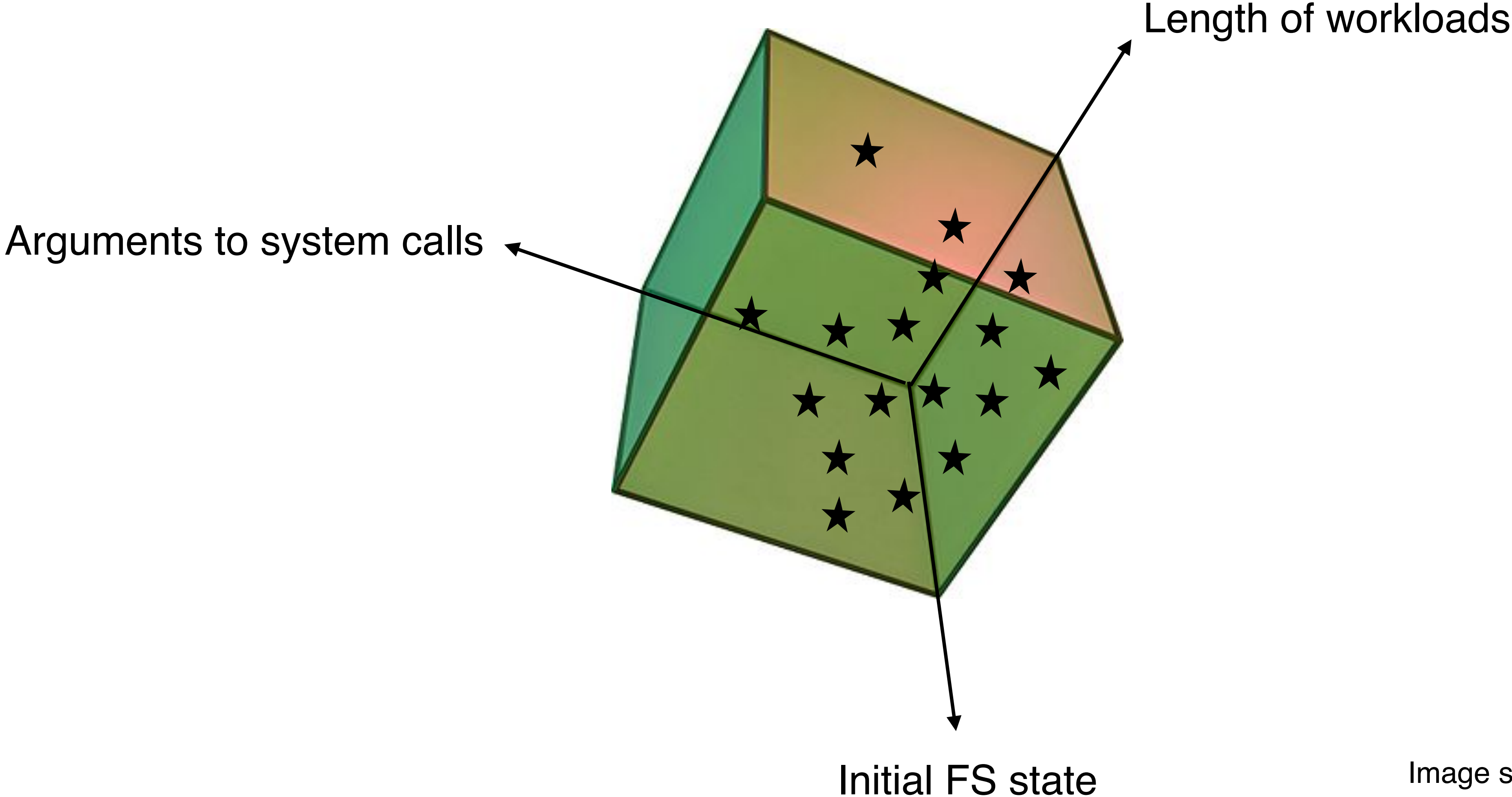
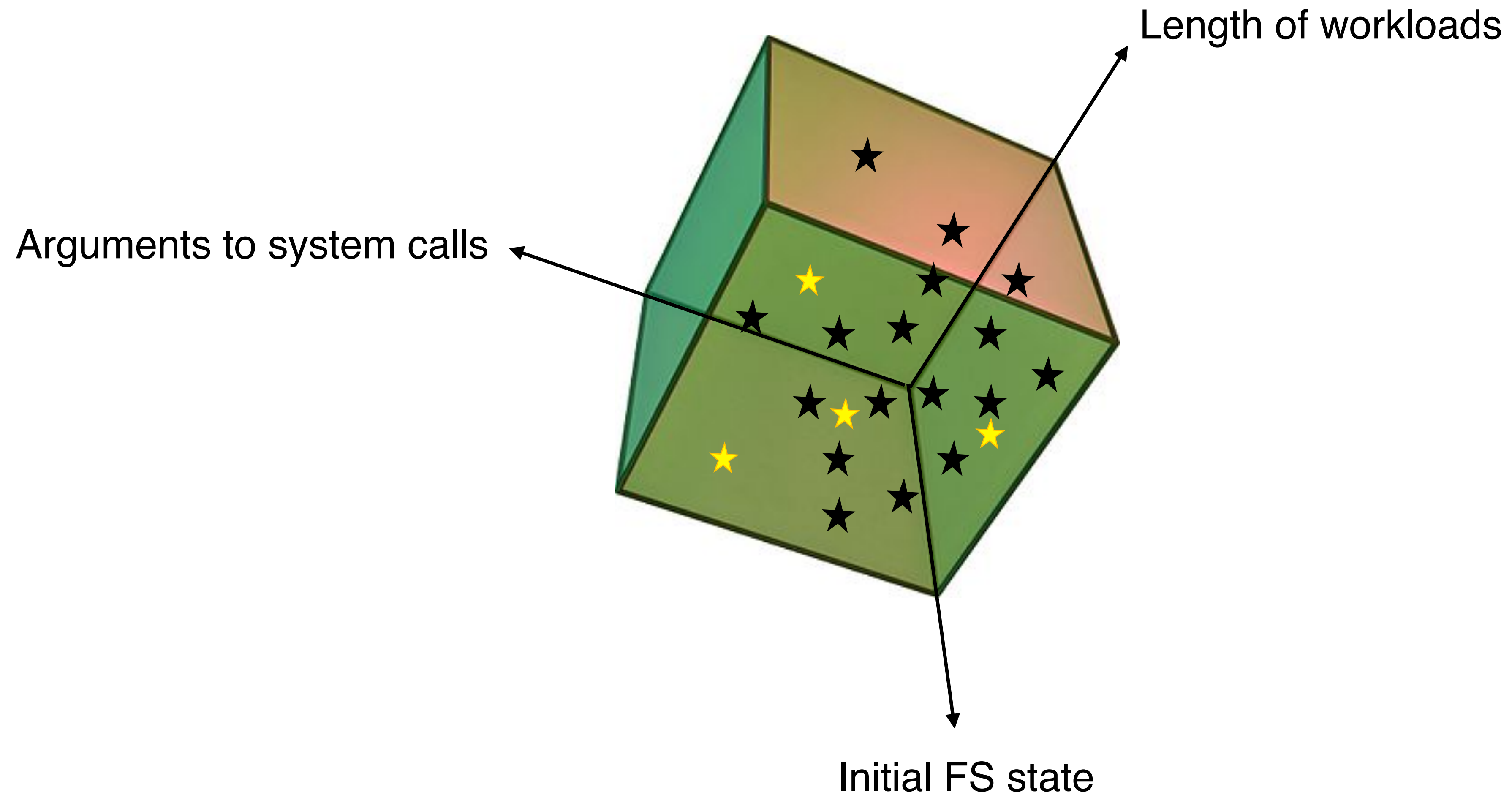
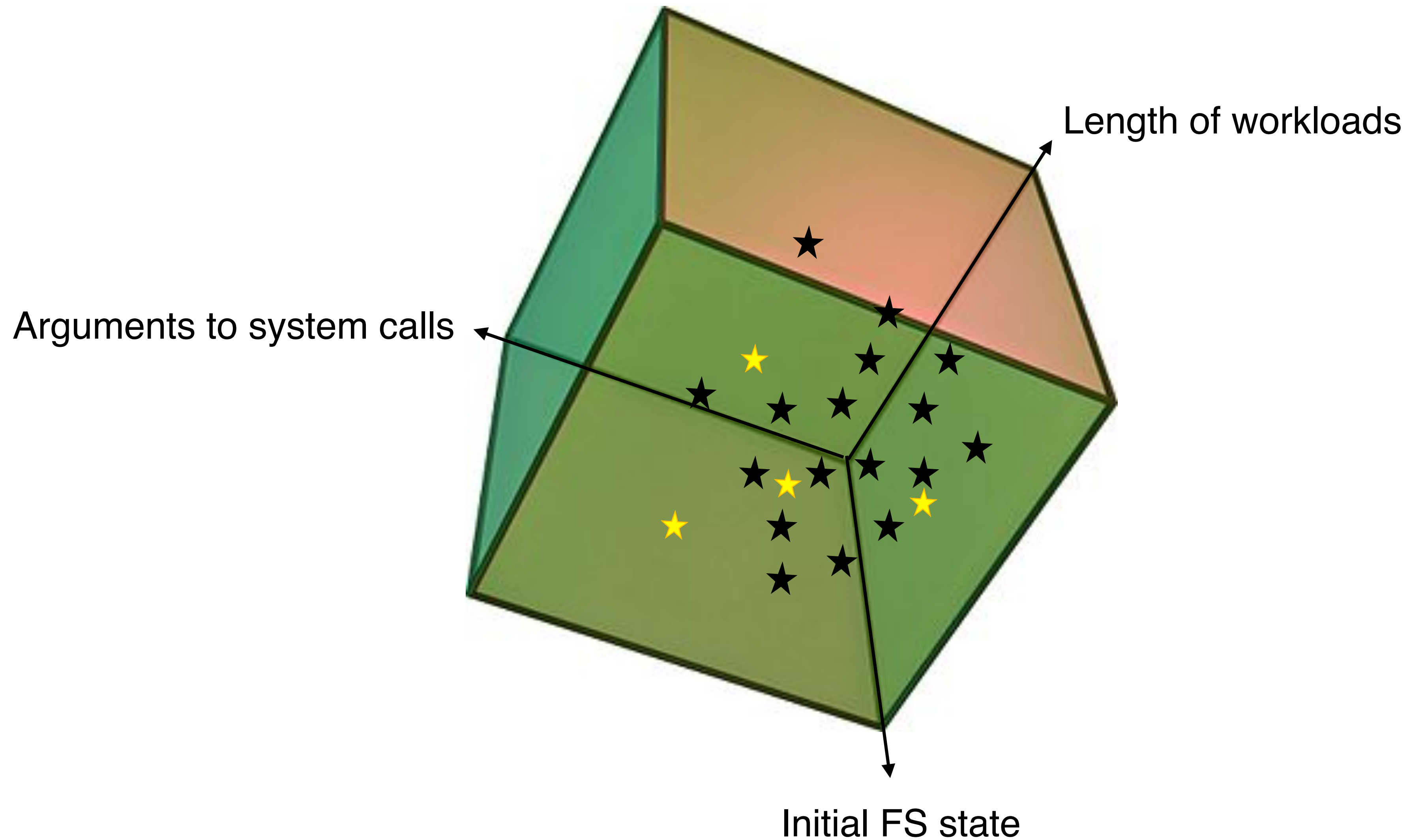


Image source: <https://en.wikipedia.org/wiki/Cube>

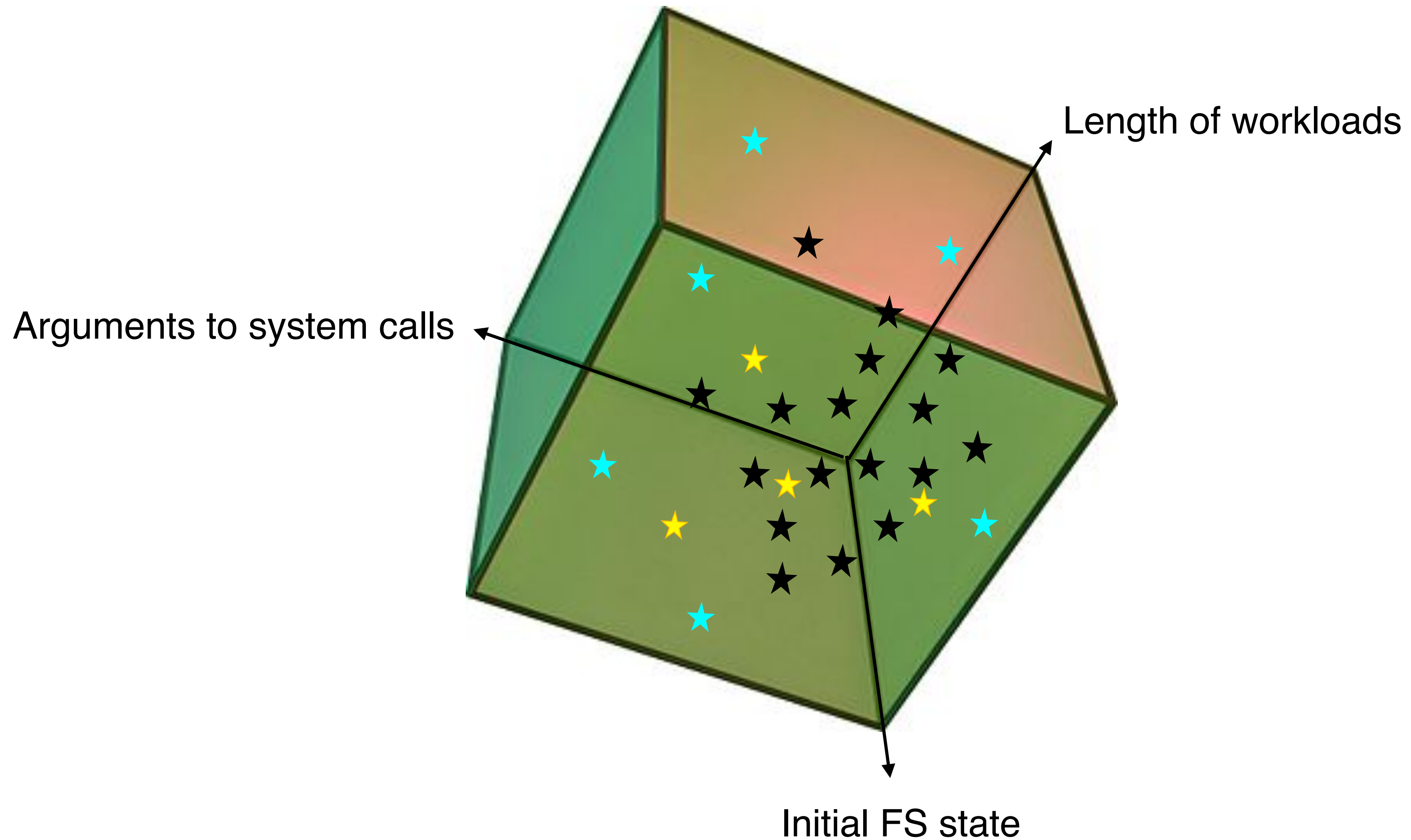
B³ : Bounded Black Box Crash Testing



B³ : Bounded Black Box Crash Testing



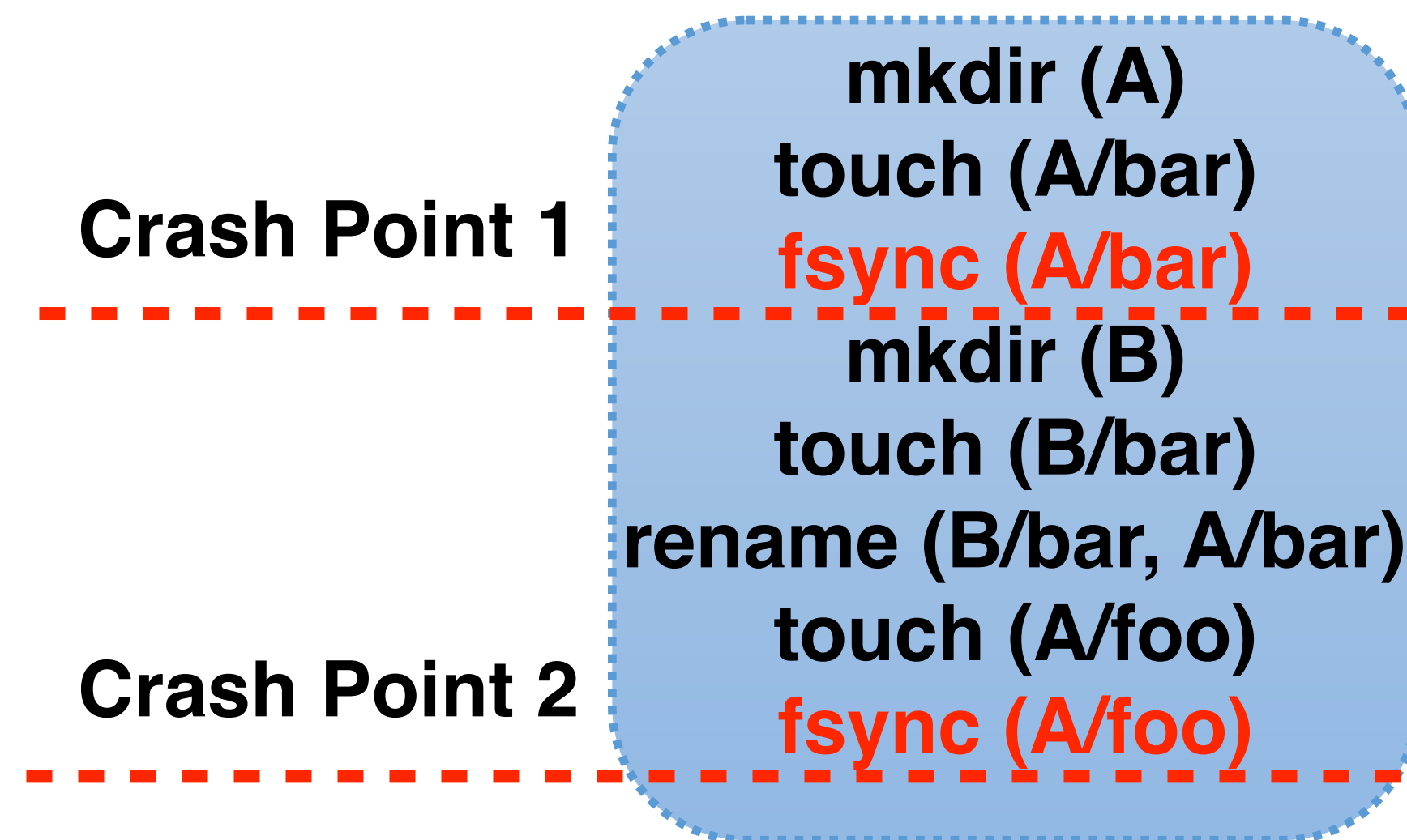
B³ : Bounded Black Box Crash Testing



B³ : Bounded Black Box Crash Testing

Choice of crash point

- Only after `fsync()`, `fdatasync()` or `sync()`
- Not in the middle of system call



- Developers are motivated to patch bugs that break semantics of persistence operations
- Crashing in the middle of system calls leads to exponentially large crash-states.

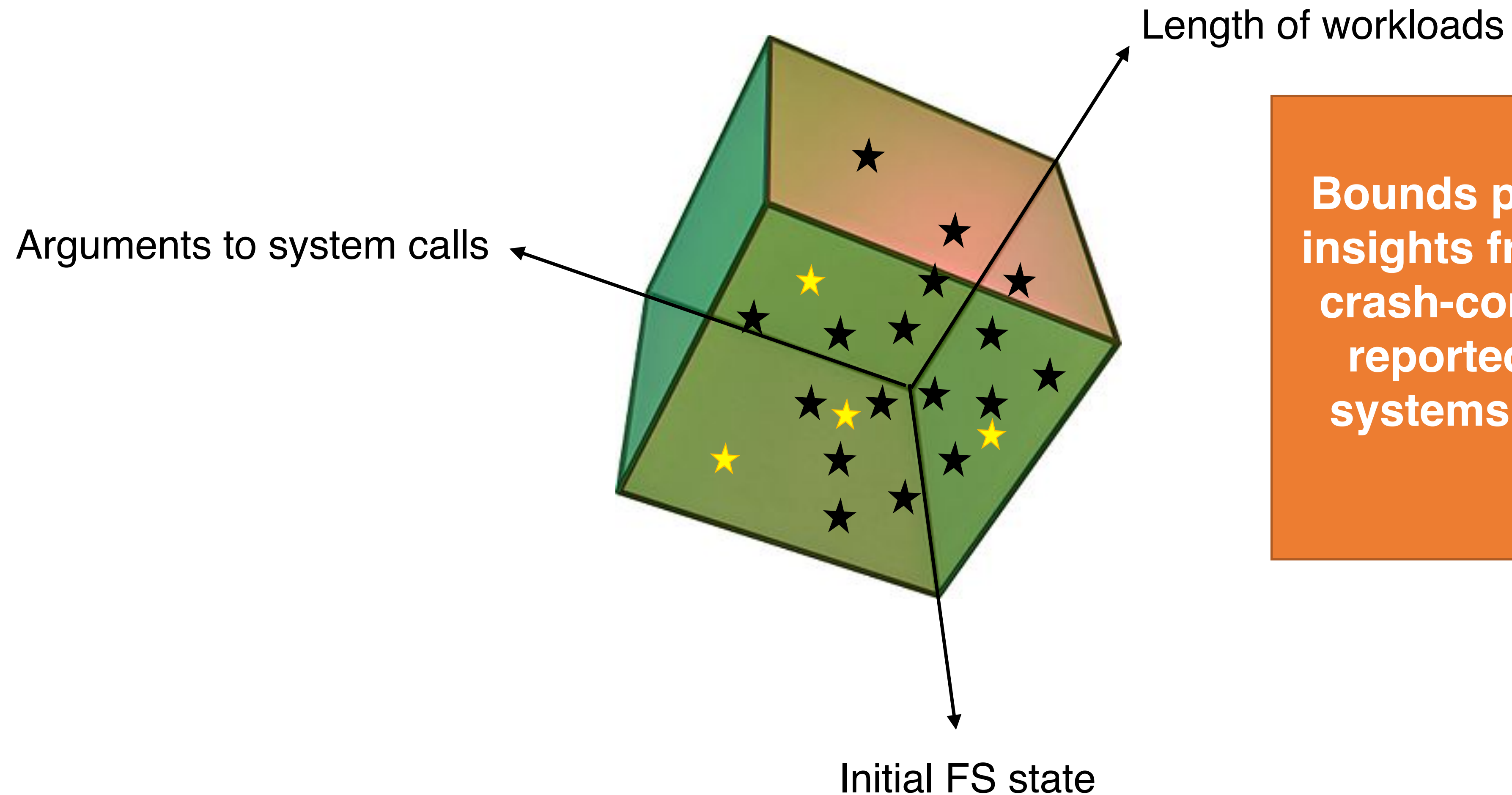
Limitations of B³

- No guarantee of finding **all** crash-consistency bugs in a filesystem
- Assumes the correct working of crash-consistency mechanism like journaling or CoW
 - Does not crash in the middle of system calls
- Can only reveal if a bug has occurred, not the reason or origin of bug.
- Needs larger compute to test higher sequence lengths

Outline

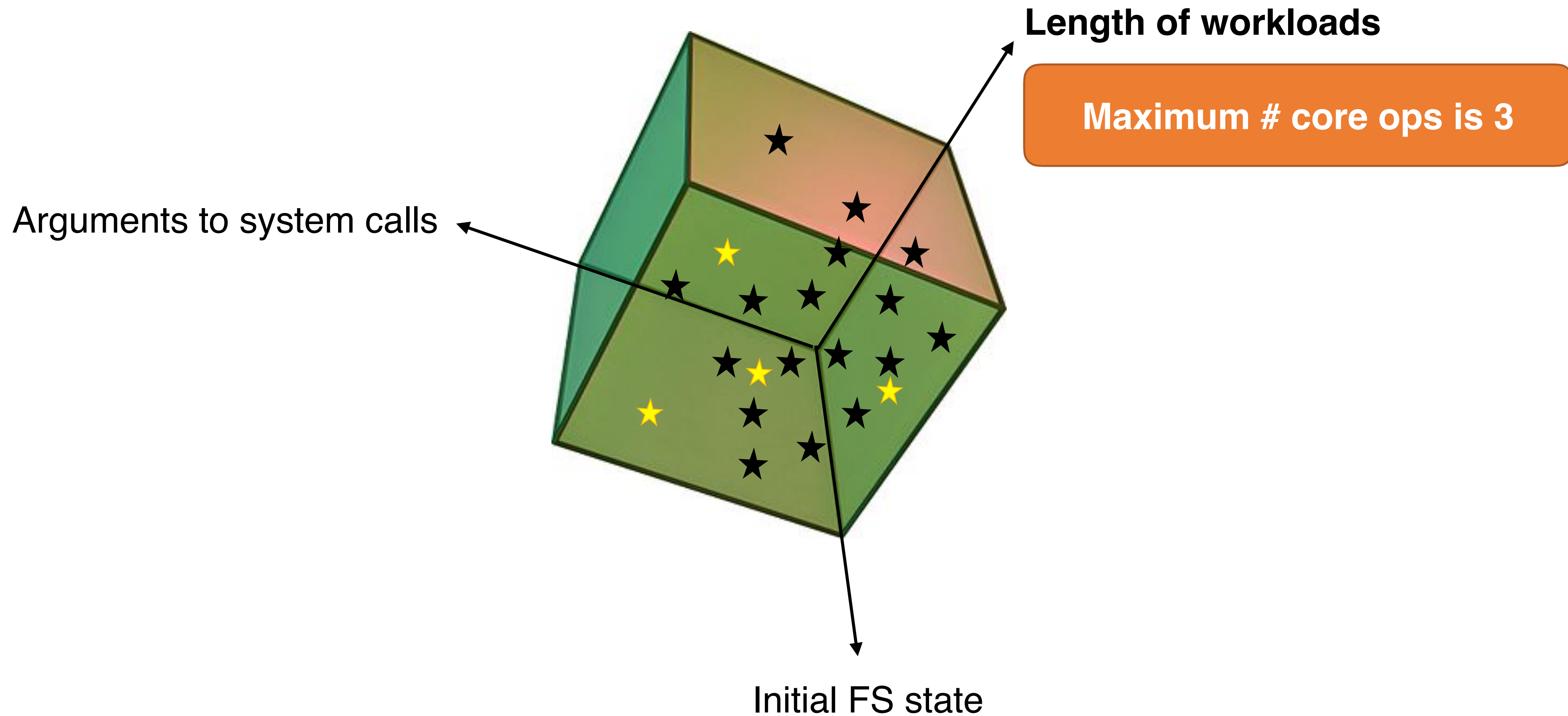
- CrashMonkey
- Bounded Black Box Crash Testing
- **Automatic Crash Explorer (ACE)**
- Demo

Bounds chosen by ACE

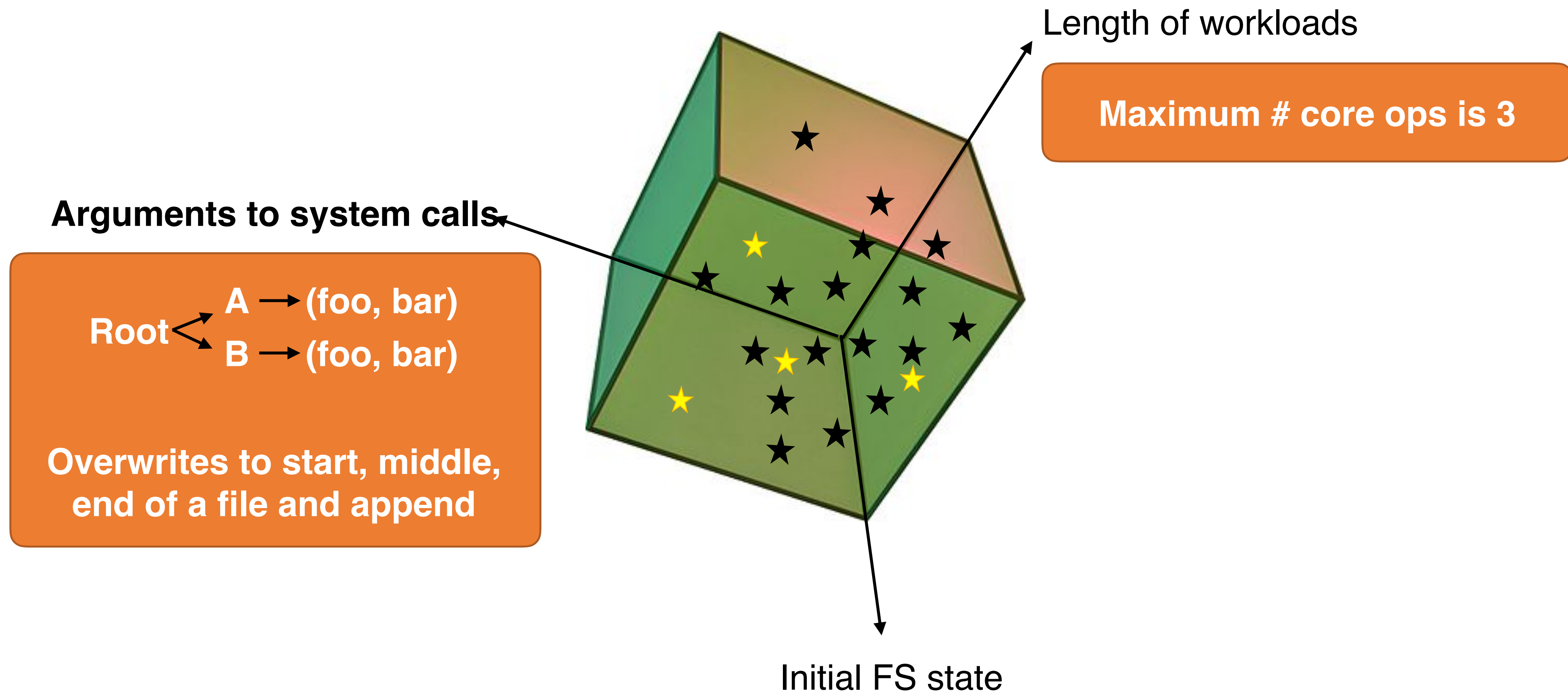


Bounds picked based on insights from the study of crash-consistency bugs reported on Linux file systems over the last 5 years

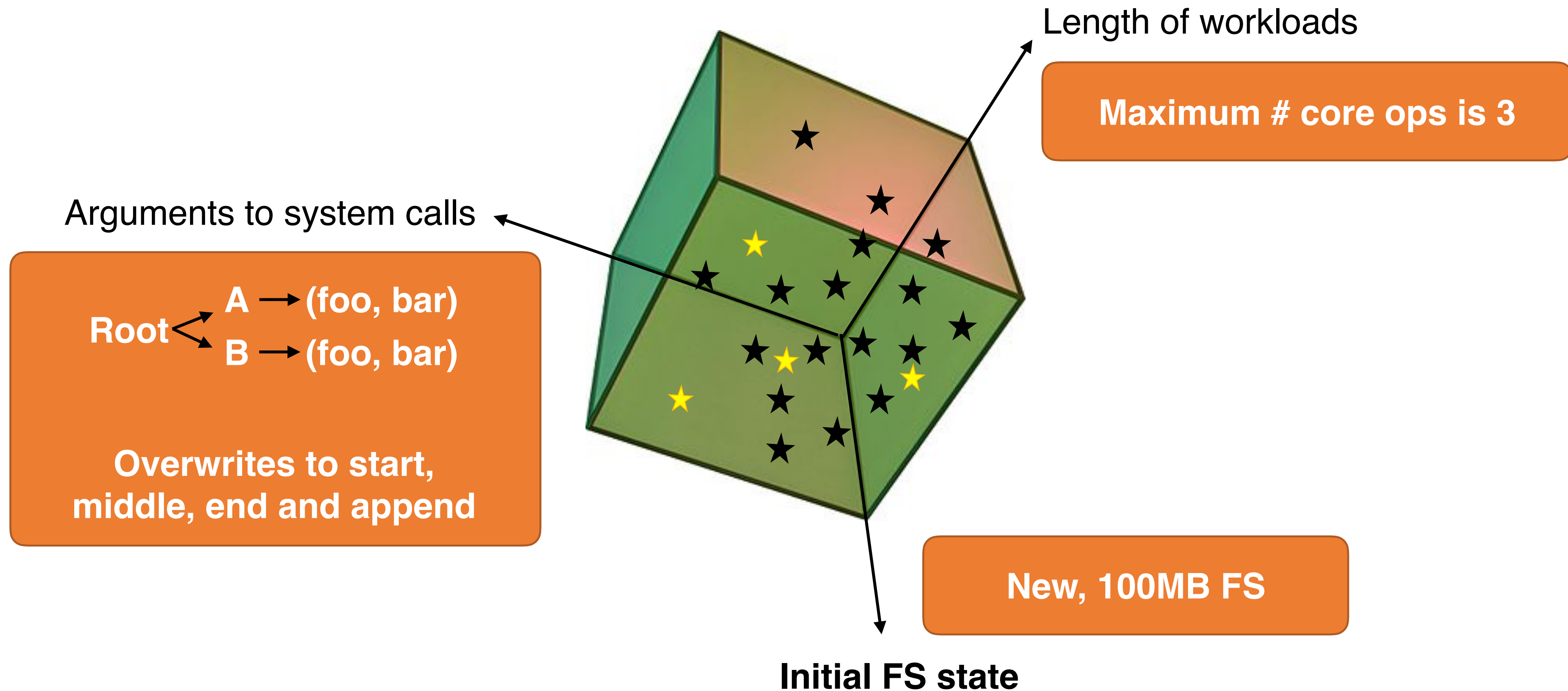
Bounds chosen by ACE



Bounds chosen by ACE



Bounds chosen by ACE

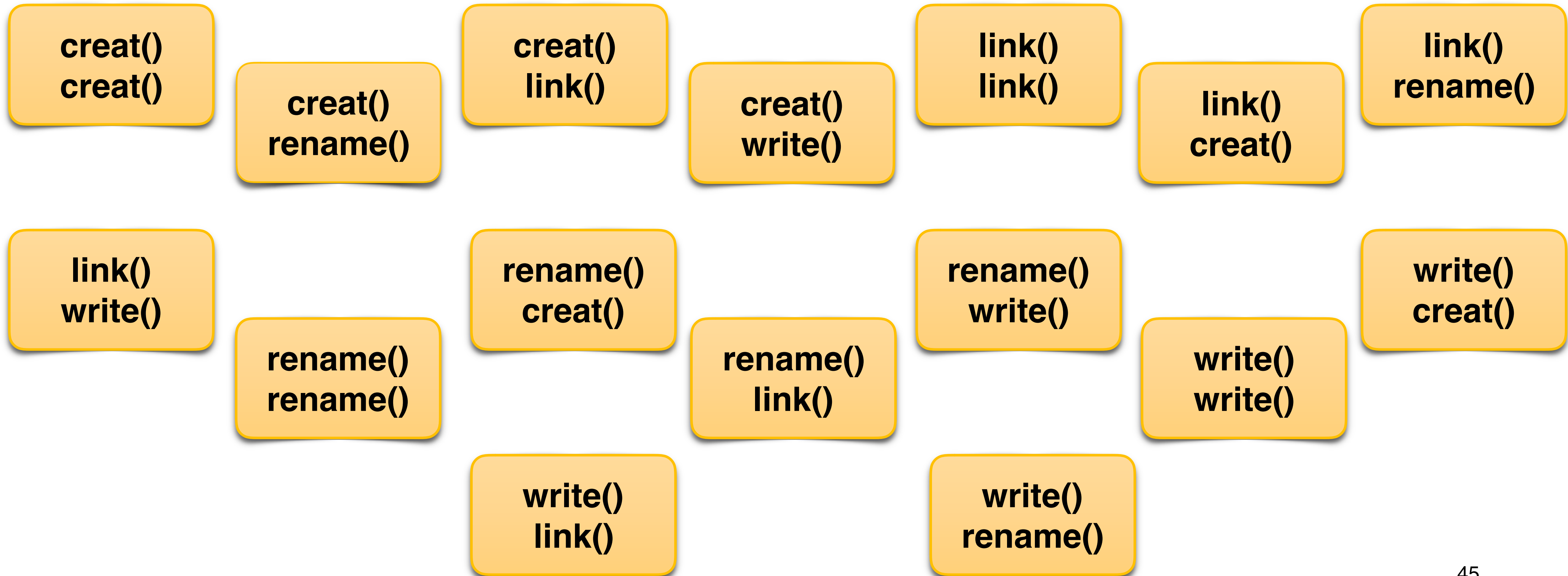


Phases of ACE

Operation Set

**creat()
link()
rename()
write()**

Generating skeletons of sequence-2. : $4*4 = 16$

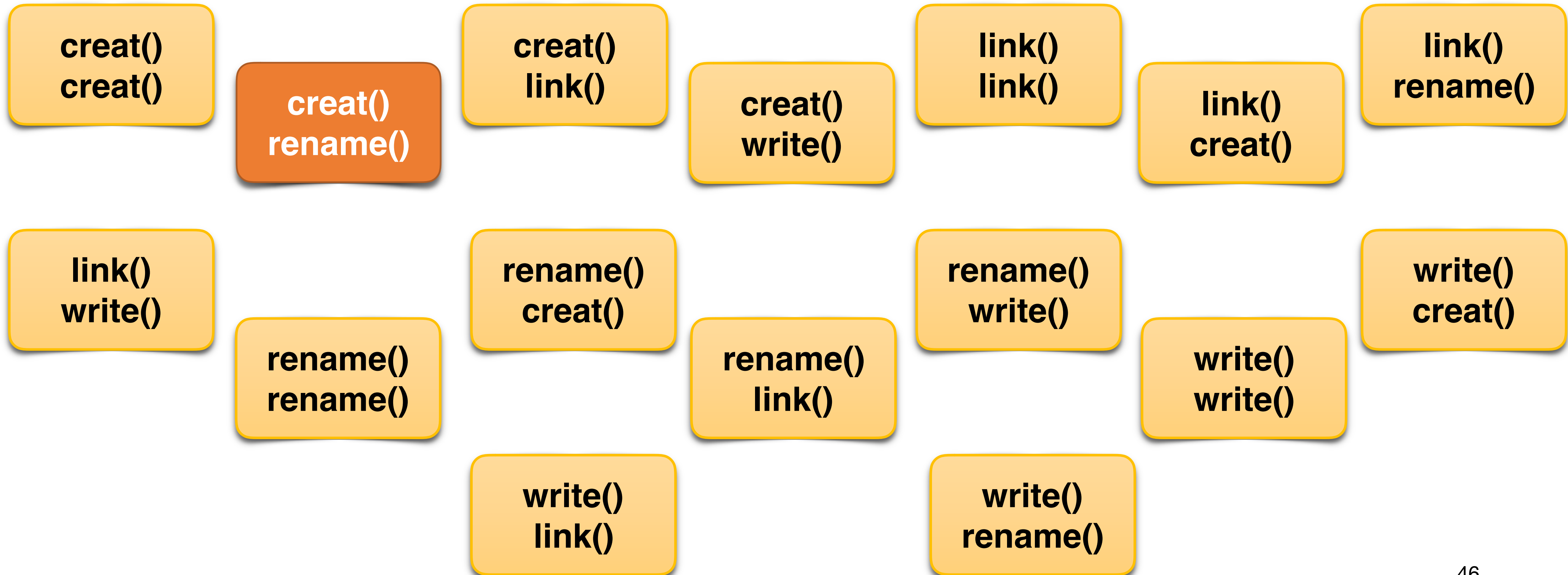


Phases of ACE

Operation Set

**creat()
link()
rename()
write()**

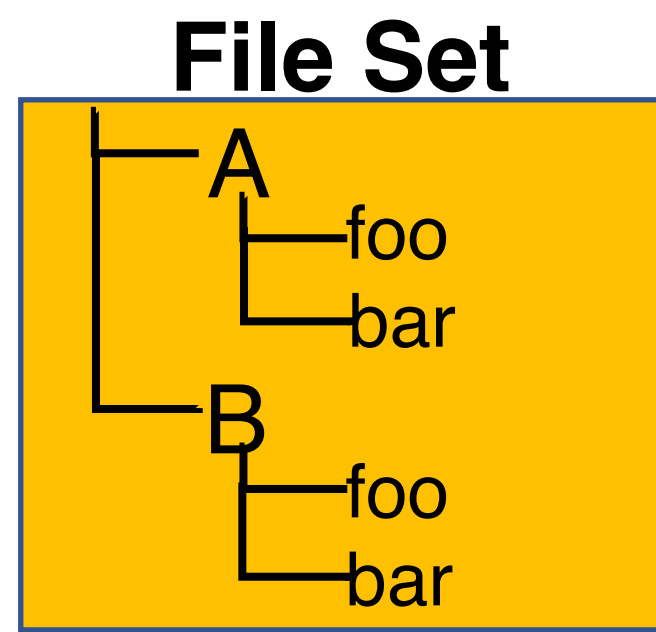
Generating skeletons of sequence-2. : $4*4 = 16$



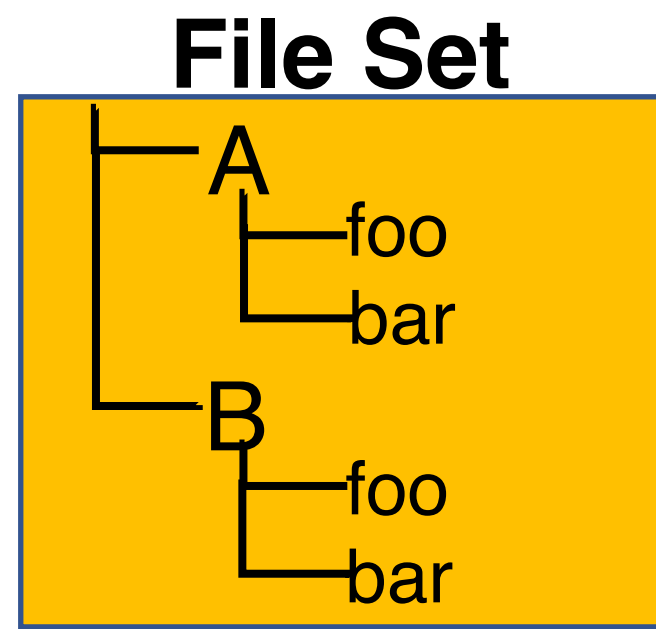
Phases of ACE

1. Select Operations

1. `creat()`
2. `rename()`



Phases of ACE



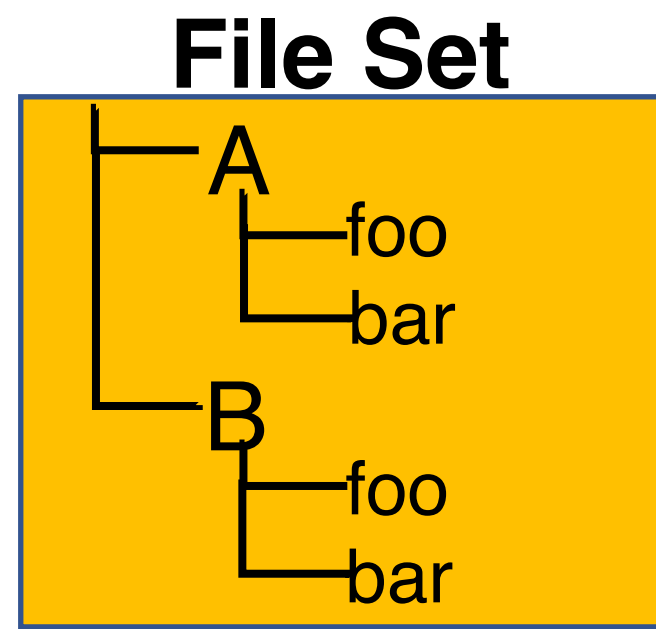
1. Select Operations

1. `creat()`
2. `rename()`

- If metadata operations, pick file or directory names
- If data operations, pick a range of offset and length

2. Select Parameters

Phases of ACE



1. Select Operations

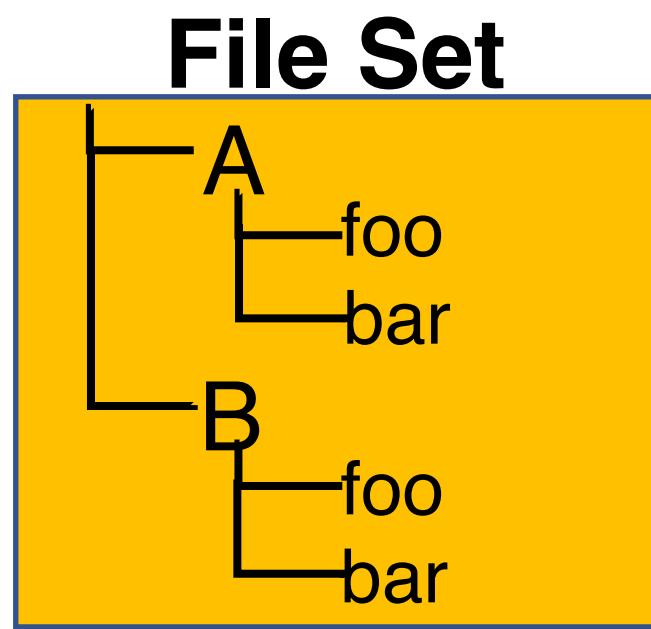
1. `creat()`
2. `rename()`

- If metadata operations, pick file or directory names
- If data operations, pick a range of offset and length

2. Select Parameters

1. `creat(A/bar)`
2. `rename(B/bar, A/bar)`

Phases of ACE



1. Select Operations

1. `creat()`
2. `rename()`

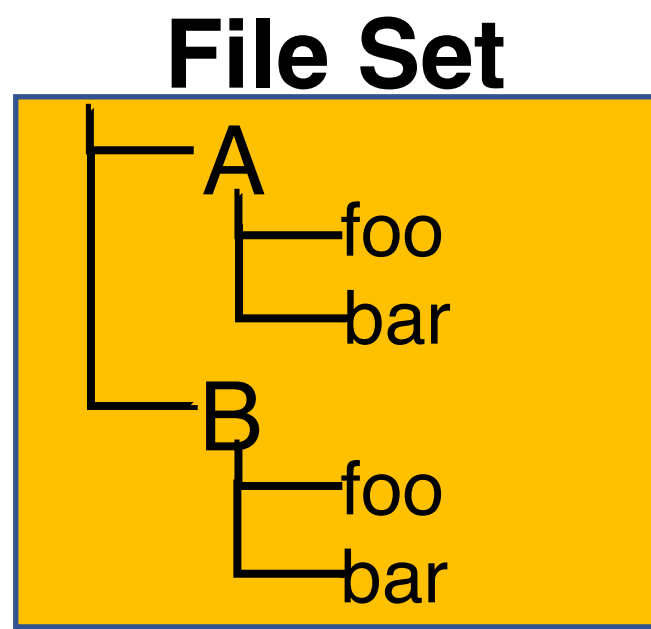
2. Select Parameters

1. `creat(A/bar)`
2. `rename(B/bar, A/bar)`

- Between each core operation, add a persistence operation
- Consistency will be checked at these points
- Parameter to the persistence function is again chosen from the file/directory pool

3. Add Persistence

Phases of ACE



1. Select Operations

1. `creat()`
2. `rename()`

2. Select Parameters

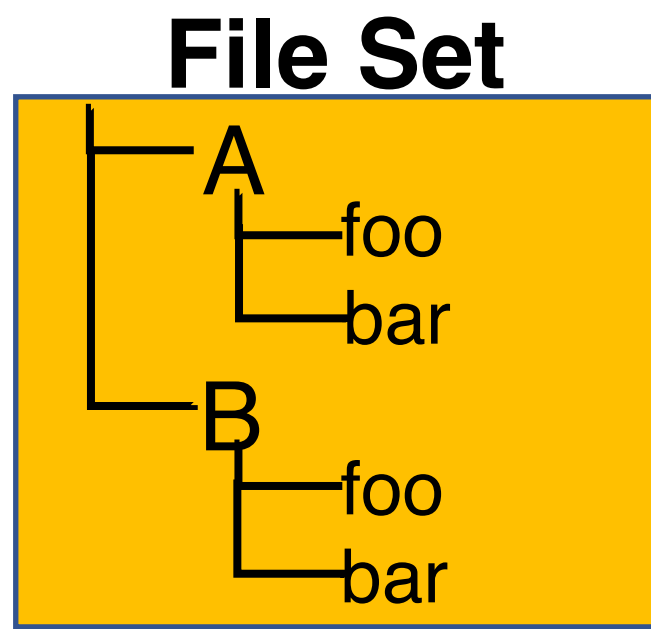
1. `creat(A/bar)`
2. `rename(B/bar, A/bar)`

- Between each core operation, add a persistence operation
- Consistency will be checked at these points
- Parameter to the persistence function is again chosen from the file/directory pool

3. Add Persistence

1. `creat(A/bar)`
`fsync(A/bar)`
2. `rename(B/bar, A/bar)`
`fsync(A/foo)`

Phases of ACE



1. Select Operations

1. `creat()`
2. `rename()`

2. Select Parameters

1. `creat(A/bar)`
2. `rename(B/bar, A/bar)`

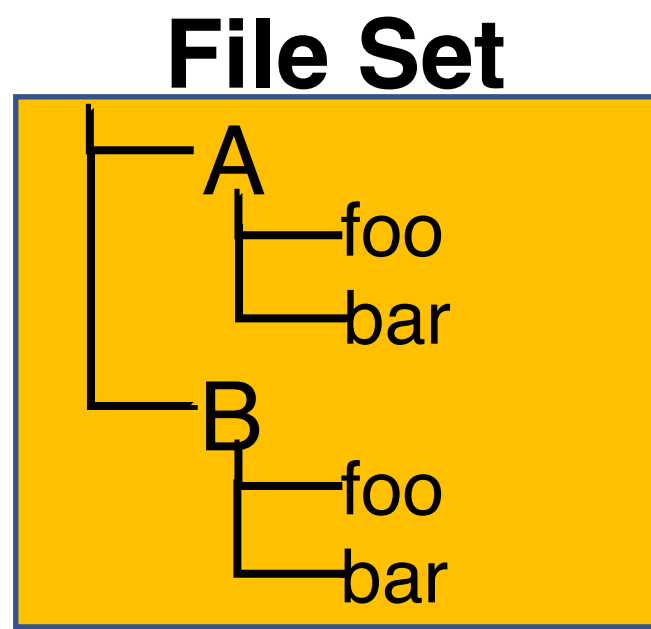
4. Add Dependencies

- Add file create/open/close to ensure the workload executes on any POSIX compliant filesystem.

3. Add Persistence

1. `creat(A/bar)`
`fsync(A/bar)`
2. `rename(B/bar, A/bar)`
`fsync(A/foo)`

Phases of ACE



1. Select Operations

1. `creat()`
2. `rename()`

2. Select Parameters

1. `creat(A/bar)`
2. `rename(B/bar, A/bar)`

4. Add Dependencies

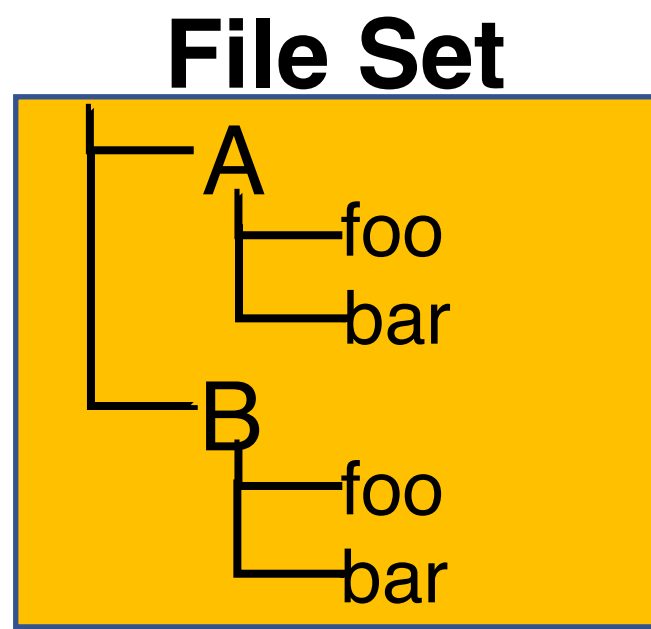
- `mkdir(A)`
1. `creat(A/bar)`
`fsync(A/bar)`
`mkdir(B)`
`creat(B/bar)`
 2. `rename(B/bar, A/bar)`
`creat(A/foo)`
`fsync(A/foo)`
`close(A/foo)`

- Add file create/open/close to ensure the workload executes on any POSIX compliant filesystem.

3. Add Persistence

1. `creat(A/bar)`
`fsync(A/bar)`
2. `rename(B/bar, A/bar)`
`fsync(A/foo)`

Phases of ACE



1. Select Operations

1. `creat()`
2. `rename()`

2. Select Parameters

1. `creat(A/bar)`
2. `rename(B/bar, A/bar)`

This workload with 2 core operations is the same workload required to trigger rename atomicity bug!

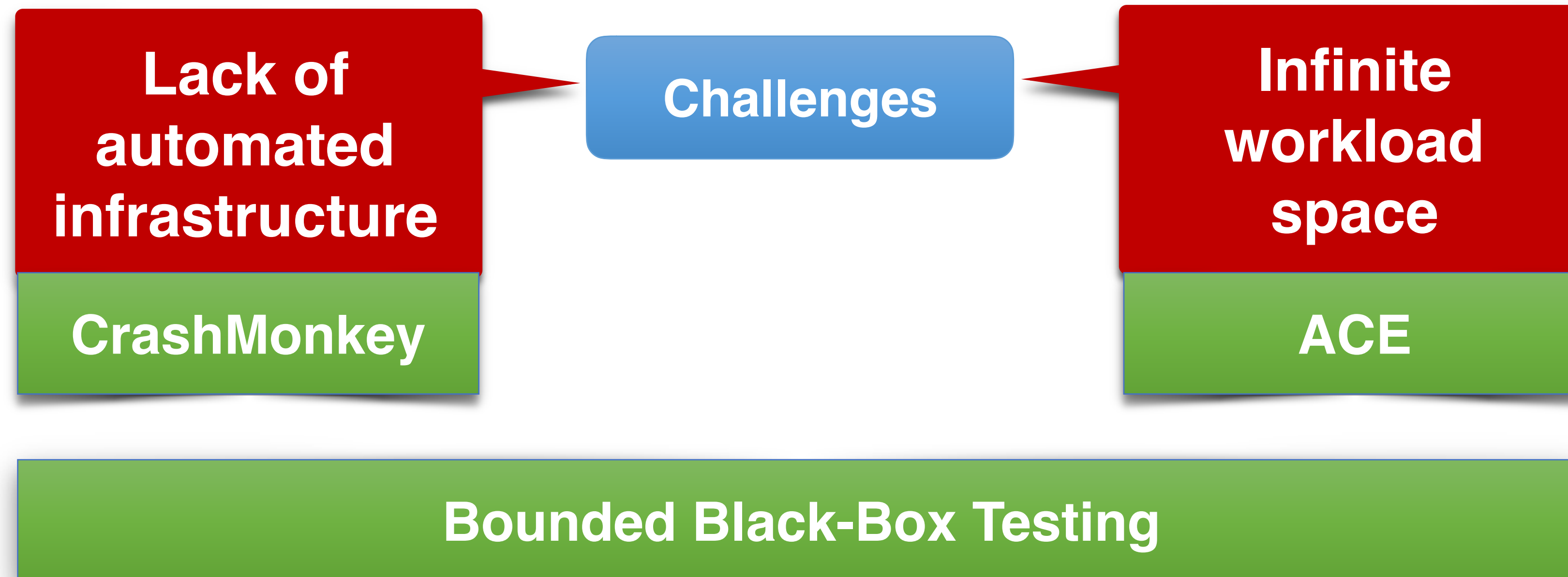
4. Add Dependencies

- `mkdir(A)`
1. `creat(A/bar)`
`fsync(A/bar)`
`mkdir(B)`
`creat(B/bar)`
 2. `rename(B/bar, A/bar)`
`creat(A/foo)`
`fsync(A/foo)`
`close(A/foo)`

3. Add Persistence

1. `creat(A/bar)`
`fsync(A/bar)`
2. `rename(B/bar, A/bar)`
`fsync(A/foo)`

Challenges with Systematic Testing



Results

- Reproduced 24/26 known bugs across ext4, btrfs and F2FS
- Found 10 new bugs across btrfs and F2FS
- Found 1 bug in a verified file system, FSCQ

Outline

- CrashMonkey
- Bounded Black Box Crash Testing
- Automatic Crash Explorer (ACE)
- **Demo**

Testing, specification, and verification



Anil Madhavapeddy @avsm · Oct 3

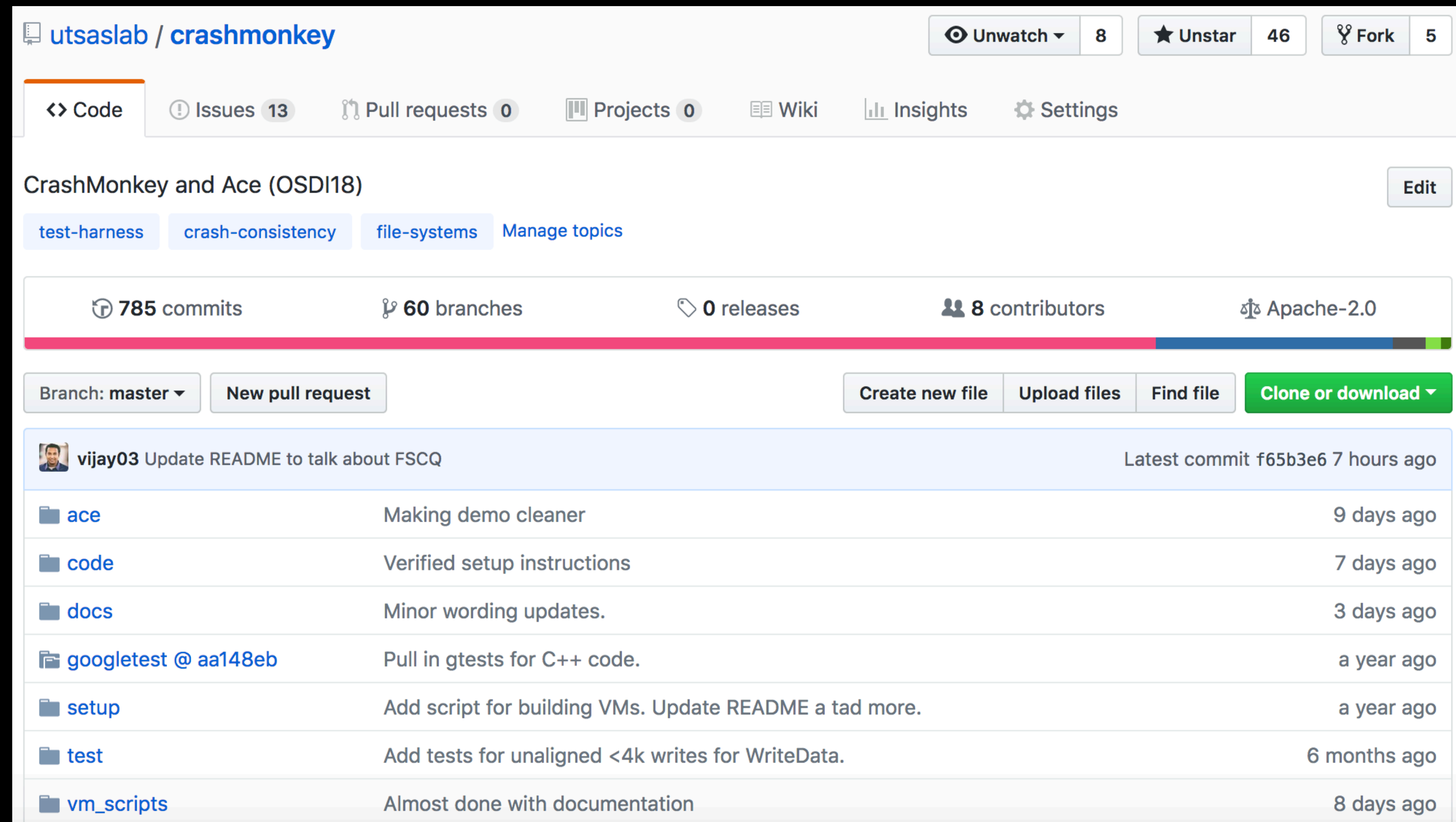


I do like how there's a systems cycle of rigorous specification, verification and testing forming — just not at the same time :-)



Bounded Black-Box Crash Testing (Poster #4)

- B³ makes exhaustive testing feasible using informed bound selection
- Easily generalizable to test larger workloads if more compute is available
- Found 10 new bugs across btrfs and F2FS, most of which existed since 2014
- Found 1 bug in FSCQ



The screenshot shows the GitHub repository page for 'utsaslab / crashmonkey'. The repository is under the 'master' branch and has 785 commits, 60 branches, 0 releases, 8 contributors, and is licensed under Apache-2.0. The repository is categorized under 'test-harness', 'crash-consistency', and 'file-systems'. The latest commit is by 'vijay03' titled 'Update README to talk about FSCQ' from 7 hours ago. The repository contains several folders: 'ace', 'code', 'docs', 'googletest @ aa148eb', 'setup', 'test', and 'vm_scripts'. Each folder has a brief description and a timestamp indicating when it was last updated.

Folder	Description	Last Updated
ace	Making demo cleaner	9 days ago
code	Verified setup instructions	7 days ago
docs	Minor wording updates.	3 days ago
googletest @ aa148eb	Pull in gtests for C++ code.	a year ago
setup	Add script for building VMs. Update README a tad more.	a year ago
test	Add tests for unaligned <4k writes for WriteData.	6 months ago
vm_scripts	Almost done with documentation	8 days ago

Thanks!
Questions?

Try our tools : <https://github.com/utsaslab/crashmonkey>

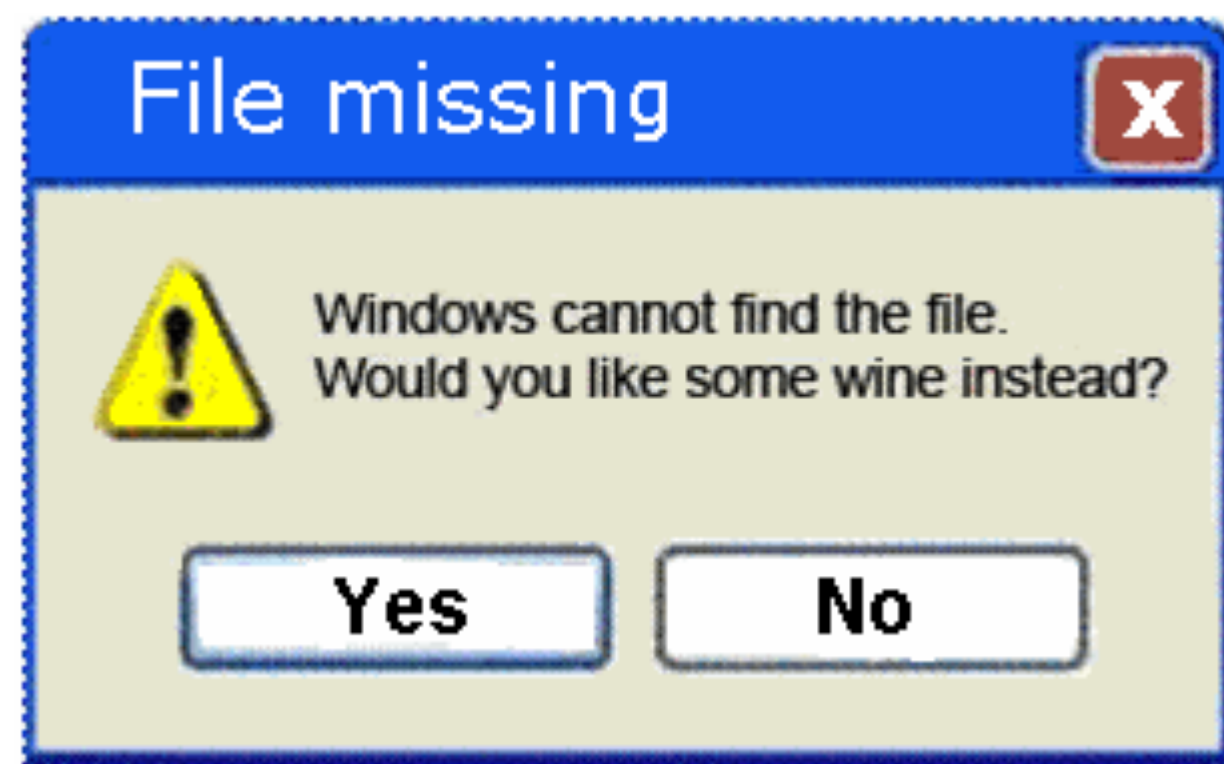
Backup slides

```
root@jayashree-VirtualBox:/home/jayashree/crashmonkey/demo/crashmonkey# ./demo.sh btrfs btrfs_out
```

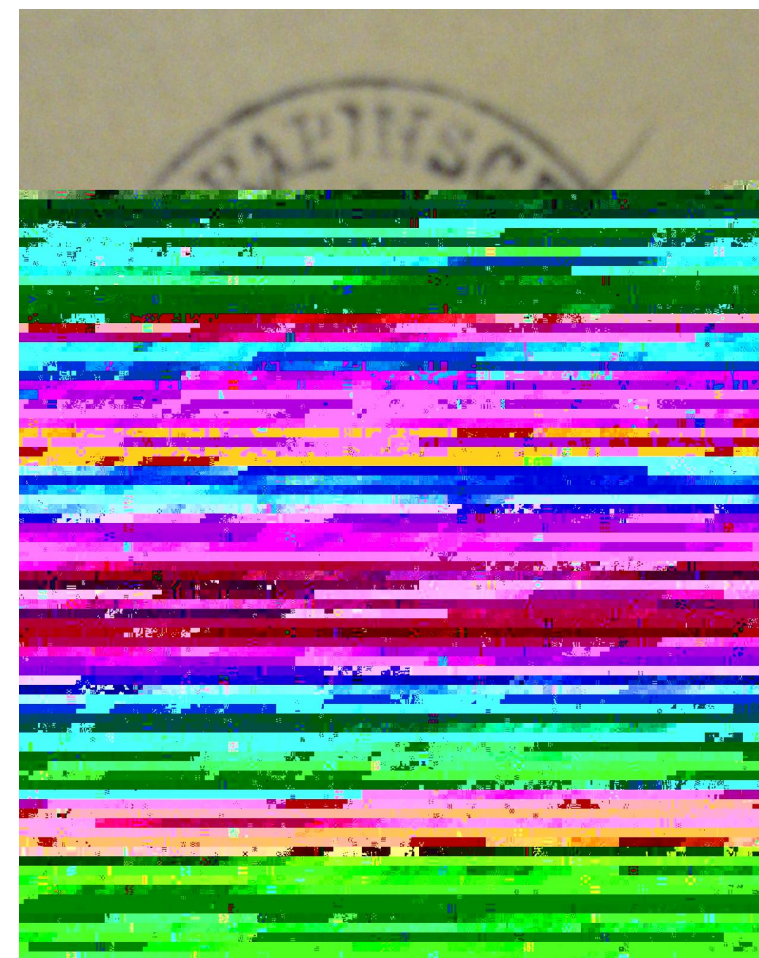
I

Crash Consistency

- Filesystem operations change multiple blocks on storage that needs to be ordered
 - Inode, bitmaps, data blocks, superblock
 - Data and metadata must be consistent on a crash



Metadata Corruption

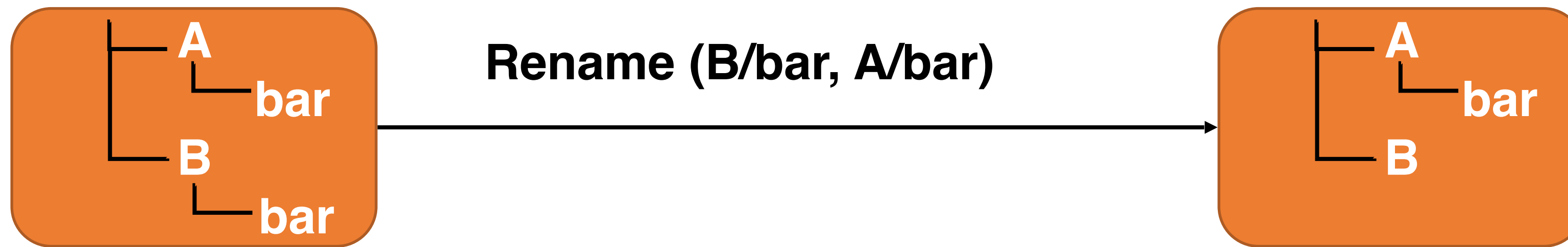


Data Corruption

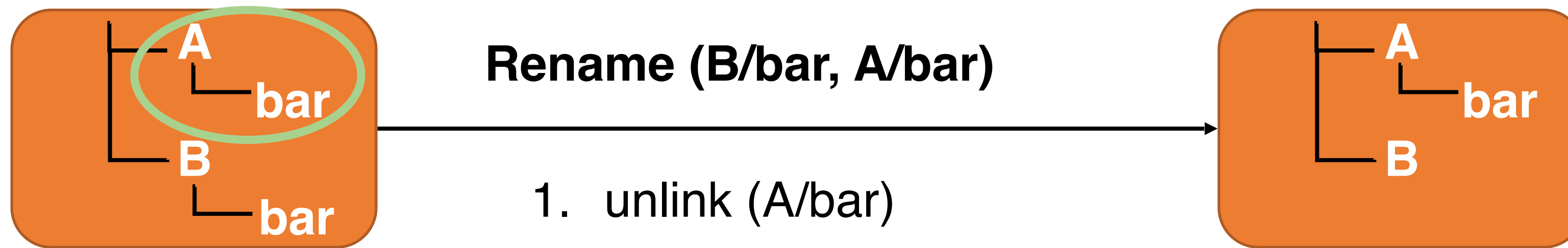


Unmountable FS

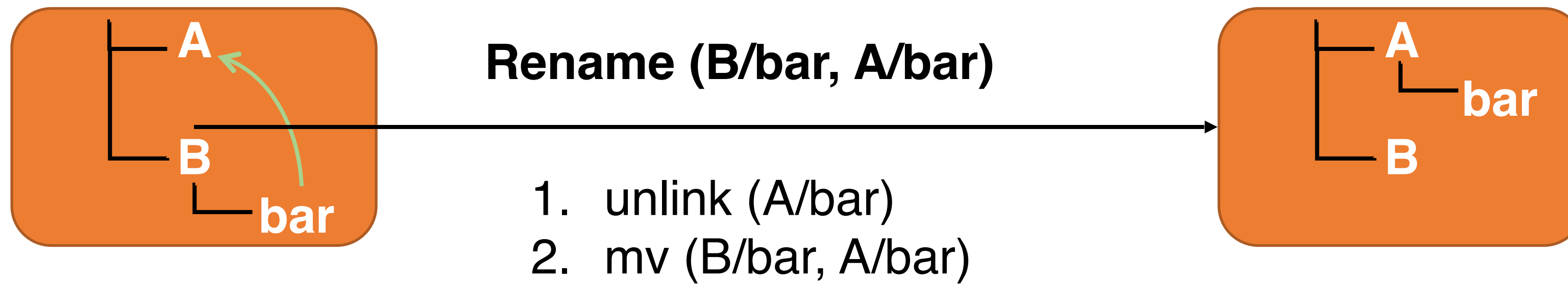
What just happened?



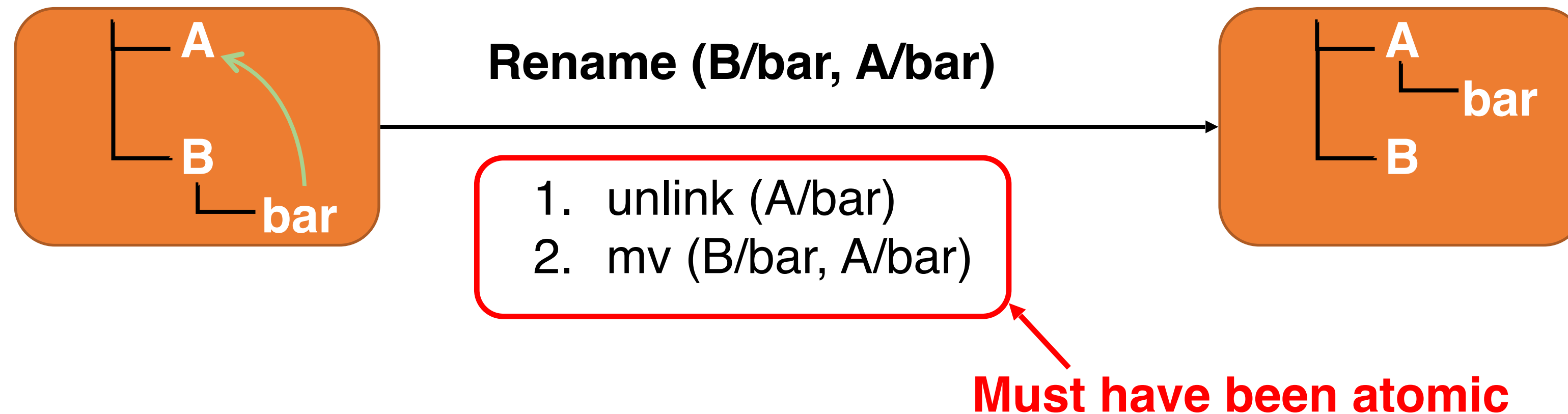
What just happened?



What just happened?



What just happened?



```
mkdir (A)
touch (A/bar)
fsync (A/bar)
mkdir (B)
touch (B/bar)
rename (B/bar, A/
bar)
touch (A/foo)
fsync (A/foo)
CRASH!
```

- fsync(A/foo) commits tx that unlinks A/bar
- Which means step 1 above is persisted, but rename is not persisted
- **End up losing file A/bar**
- Exists in the kernel since 2014

Study of crash consistency bugs in the wild

- Study the workload pattern and impacts of crash consistency bugs reported in the **past 5 years**
 - Kernel mailing lists
 - Crash consistency tests submitted to xfstests
- **26 unique bugs** across ext4, F2FS, and btrfs

Study of crash consistency bugs in the wild

Consequence	# bugs
Corruption	17
Data inconsistency	6
Unmountable FS	3
Total	26

Filesystem	# bugs
Ext4	2
F2FS	2
btrfs	24
Total	28

# ops	# bugs
1	3
2	14
3	9
Total	26

Study of crash consistency bugs in the wild

Consequence	# bugs	Filesystem	# bugs	# ops	# bugs
Corruption	17	Ext4	2	1	3
Data inconsistency	6	F2FS	2	2	14
Unmountable FS	3	btrfs	24	3	9
Total	26	Total	28	Total	26

1. Crash consistency bugs are hard to find

- Bugs have been around in the kernel for up to 7 years before being identified and patched
- Usually involve reuse of files/ directories

Study of crash consistency bugs in the wild

Consequence	# bugs	Filesystem	# bugs	# ops	# bugs
Corruption	17	Ext4	2	1	3
Data inconsistency	6	F2FS	2	2	14
Unmountable FS	3	btrfs	24	3	9
Total	26	Total	28	Total	26

1. Crash consistency bugs are hard to find
- 2. Small workloads are sufficient to reveal bugs**
 - 2-3 core operations on a new, empty file-system

Study of crash consistency bugs in the wild

Consequence	# bugs	Filesystem	# bugs	# ops	# bugs
Corruption	17	Ext4	2	1	3
Data inconsistency	6	F2FS	2	2	14
Unmountable FS	3	btrfs	24	3	9
Total	26	Total	28	Total	26

1. Crash consistency bugs are hard to find
2. Small workloads are sufficient to reveal bugs
3. **Crash after persistence points**
 - Sufficient to crash after a call to `fsync()`, `fdatasync()`, or `sync()`

Study of crash consistency bugs in the wild

Consequence	# bugs	Filesystem	# bugs	# ops	# bugs
Corruption	17	Ext4	2	1	3
Data inconsistency	6	F2FS	2	2	14
Unmountable FS	3	btrfs	24	3	9
Total	26	Total	28	Total	26

1. Crash consistency bugs are hard to find
2. Small workloads are sufficient to reveal bugs
3. Crash after persistence points
4. **Systematic testing is required**

Study of crash consistency bugs in the wild

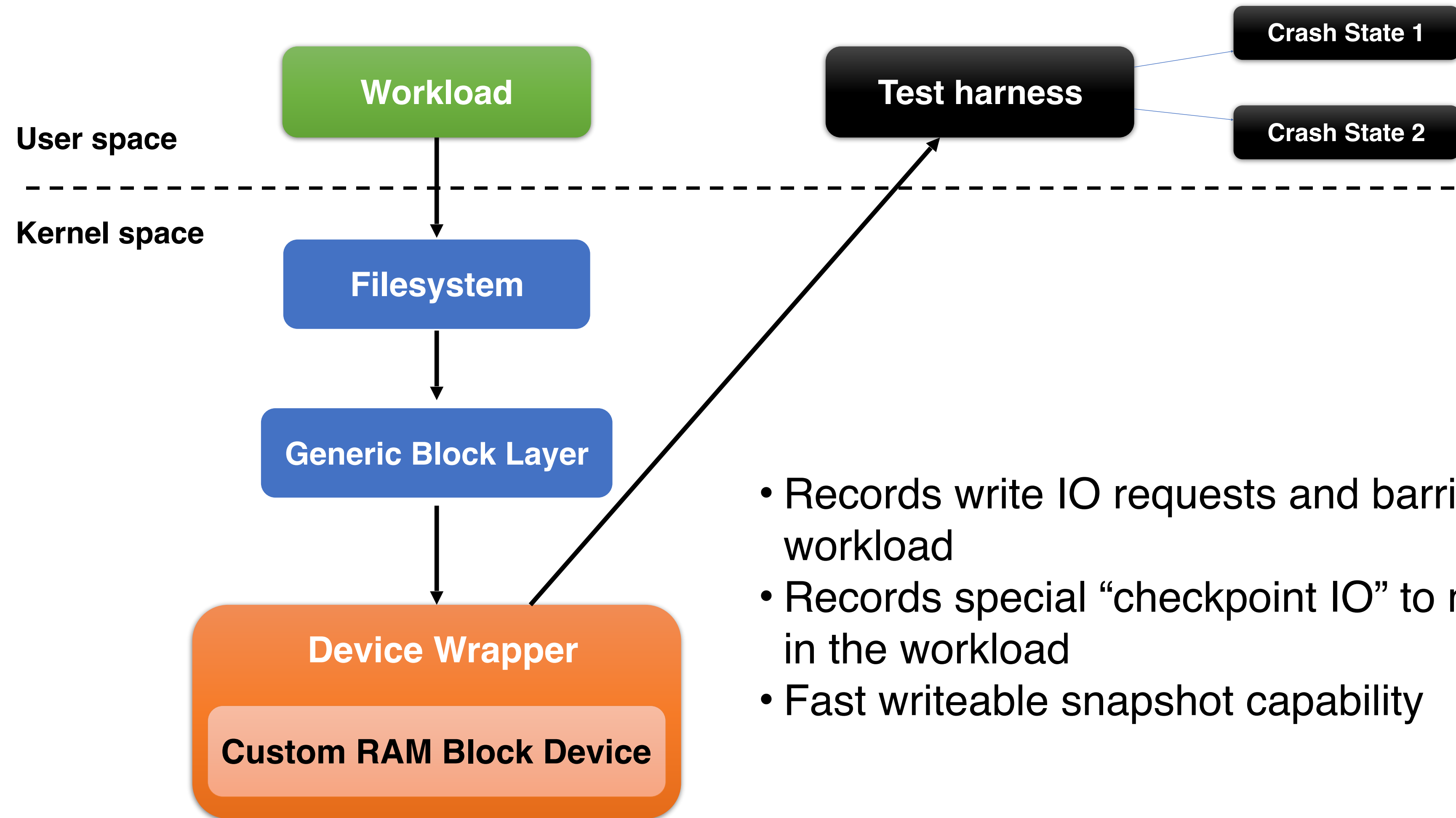
Consequence	# bugs	Filesystem	# bugs	# ops	# bugs
Corruption	17	Ext4	2	1	3
Data inconsistency	6	F2FS	2	2	14
Unmountable FS	3	btrfs	24	3	9
Total	26	Total	28	Total	26

1. Crash consistency bugs are hard to find
2. Small workloads are sufficient to reveal bugs
3. Crash after persistence points
4. **Systematic testing is required**

Fallocate : punch_hole : 2015

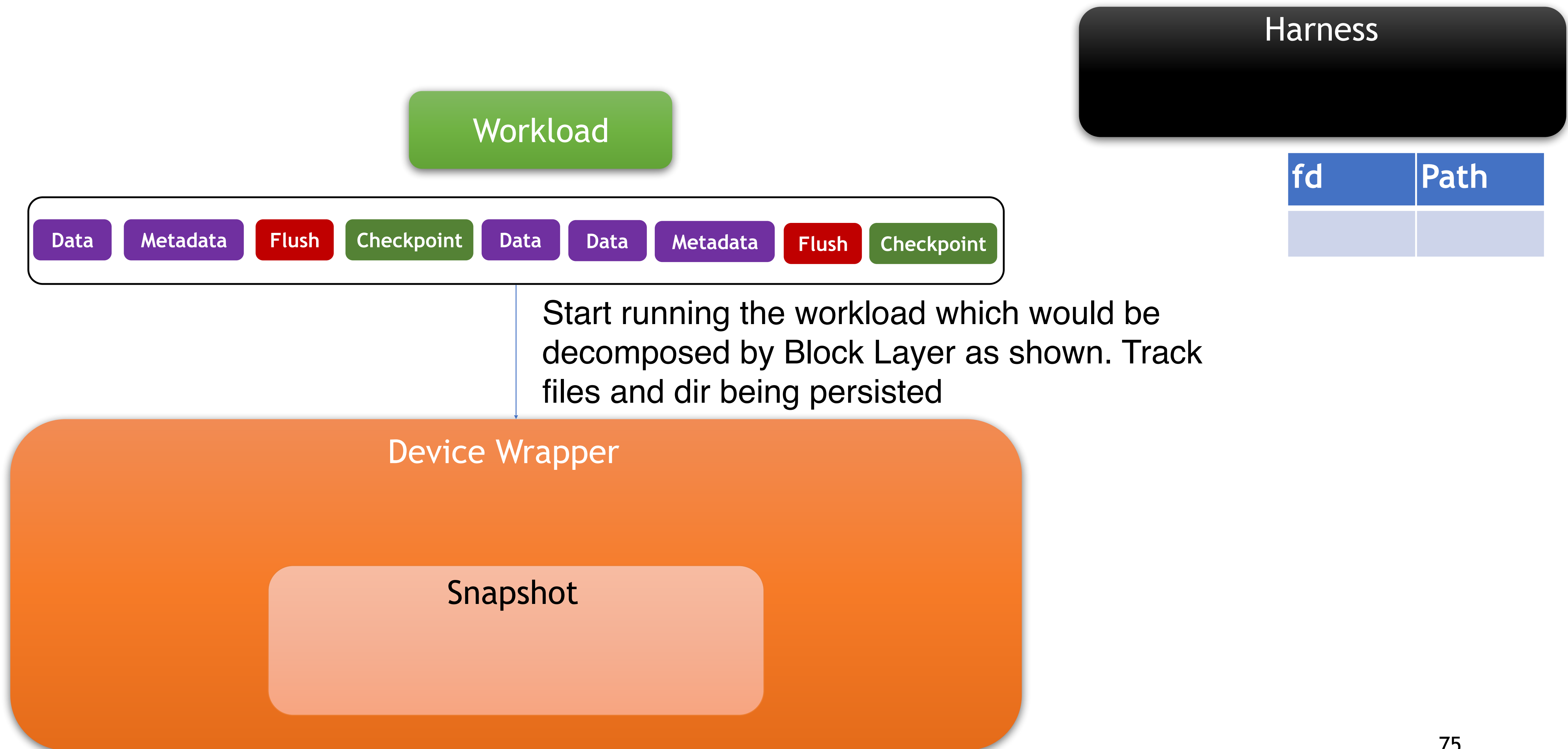
Fallocate : zero_range : 2018

CrashMonkey Internals

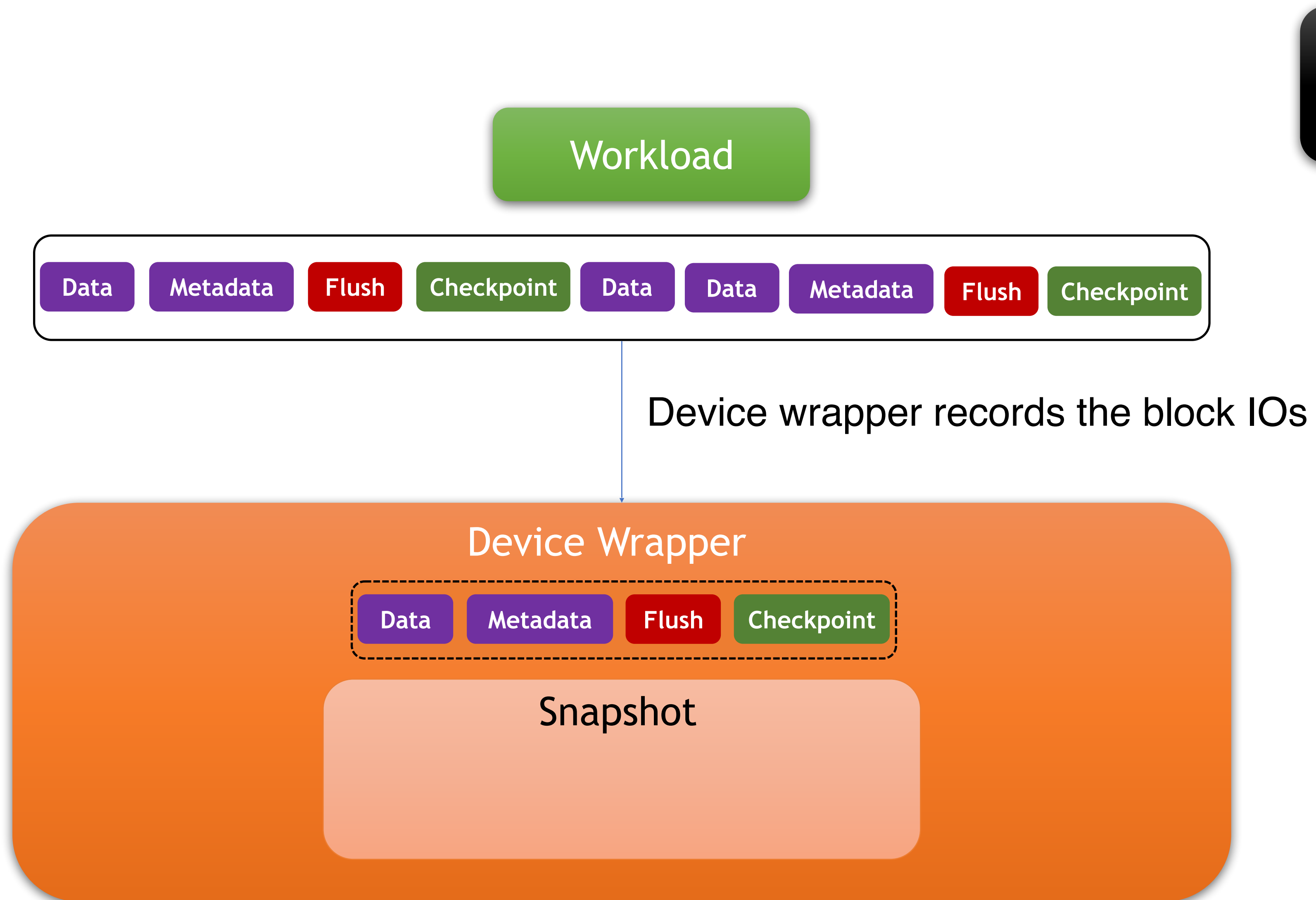


- Records write IO requests and barriers (flush/FUA) in the workload
- Records special “checkpoint IO” to mark persistence points in the workload
- Fast writeable snapshot capability

CrashMonkey in Action



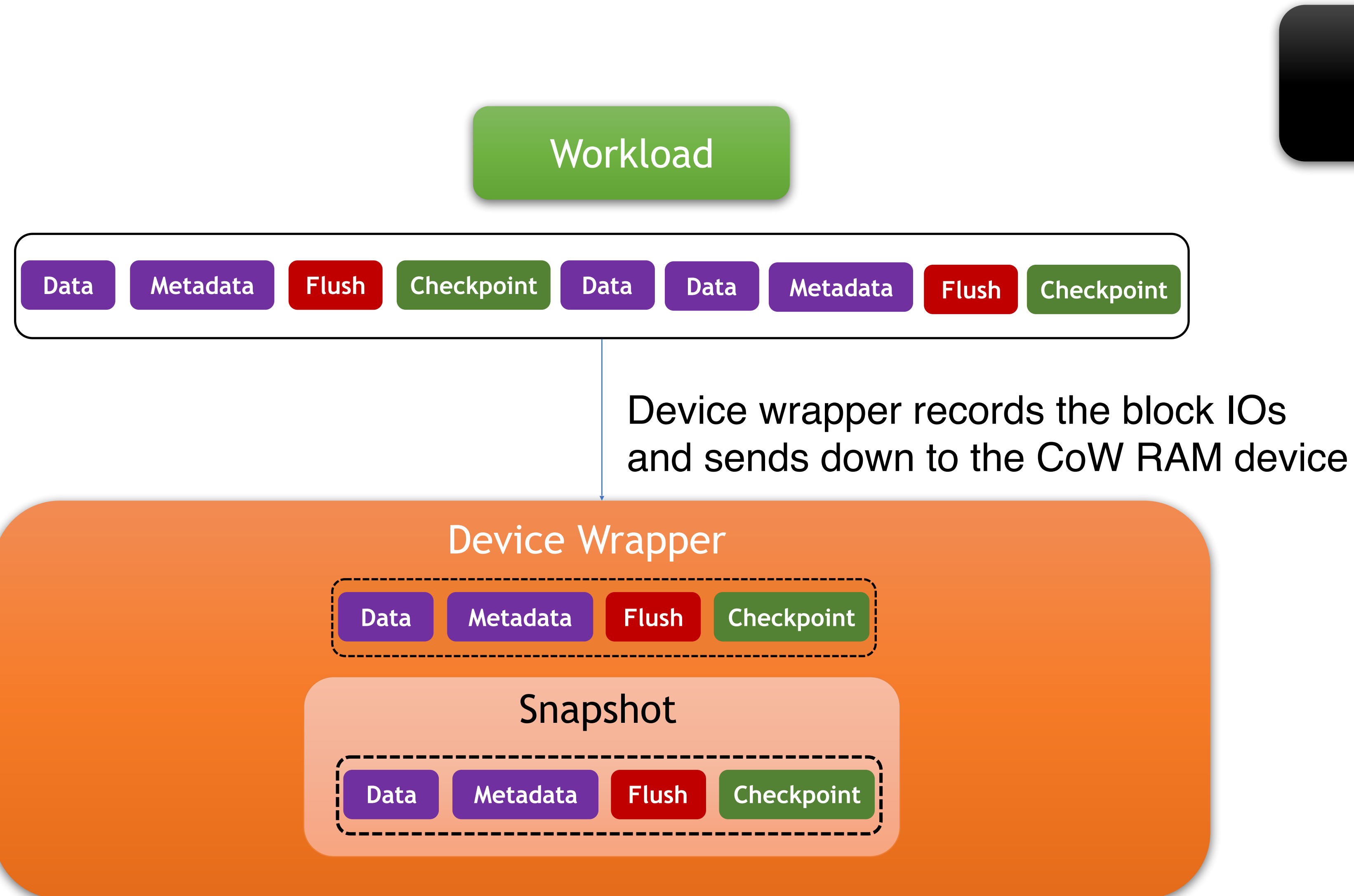
CrashMonkey in Action : Profiling



Harness

fd	Path
13	/a/b

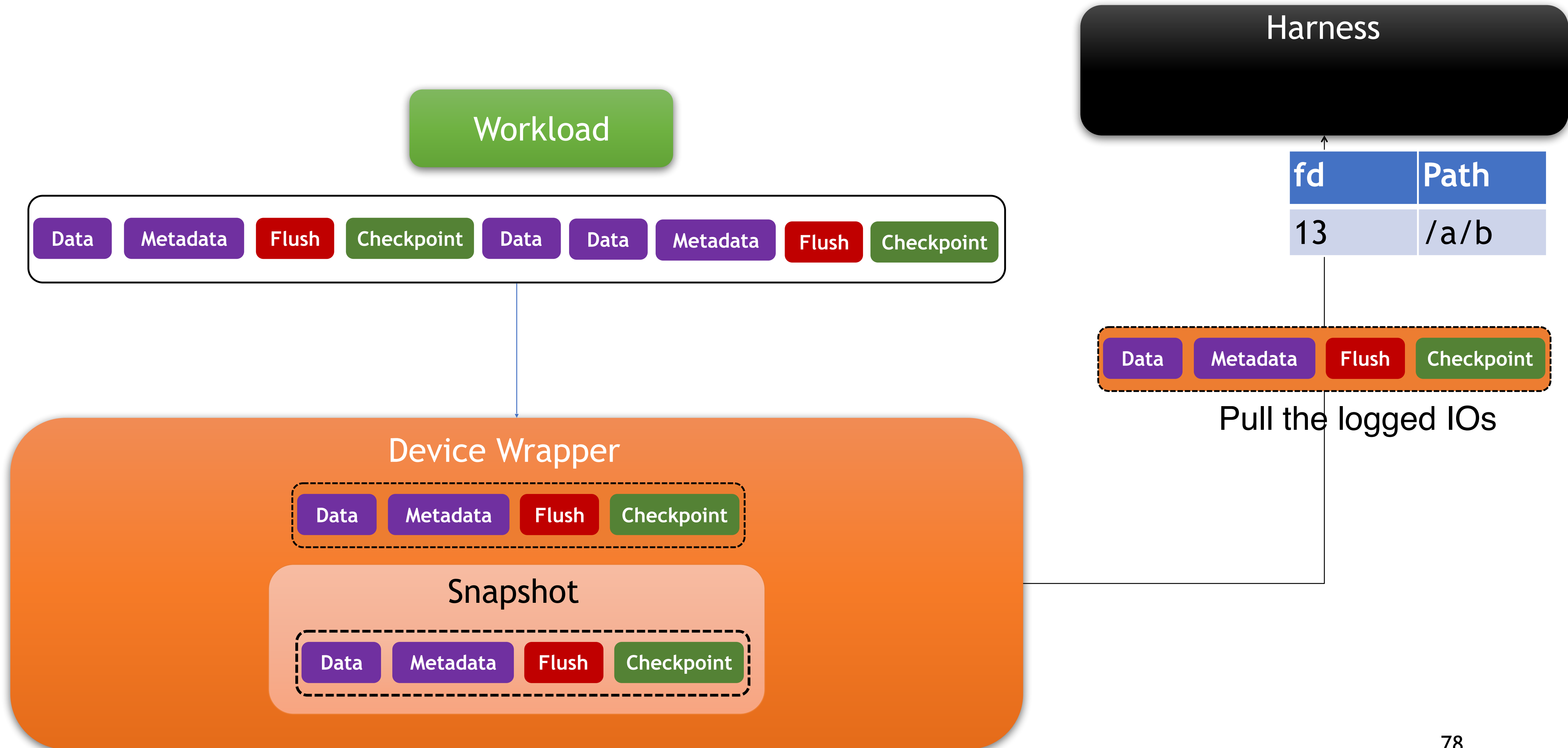
CrashMonkey in Action : Profiling



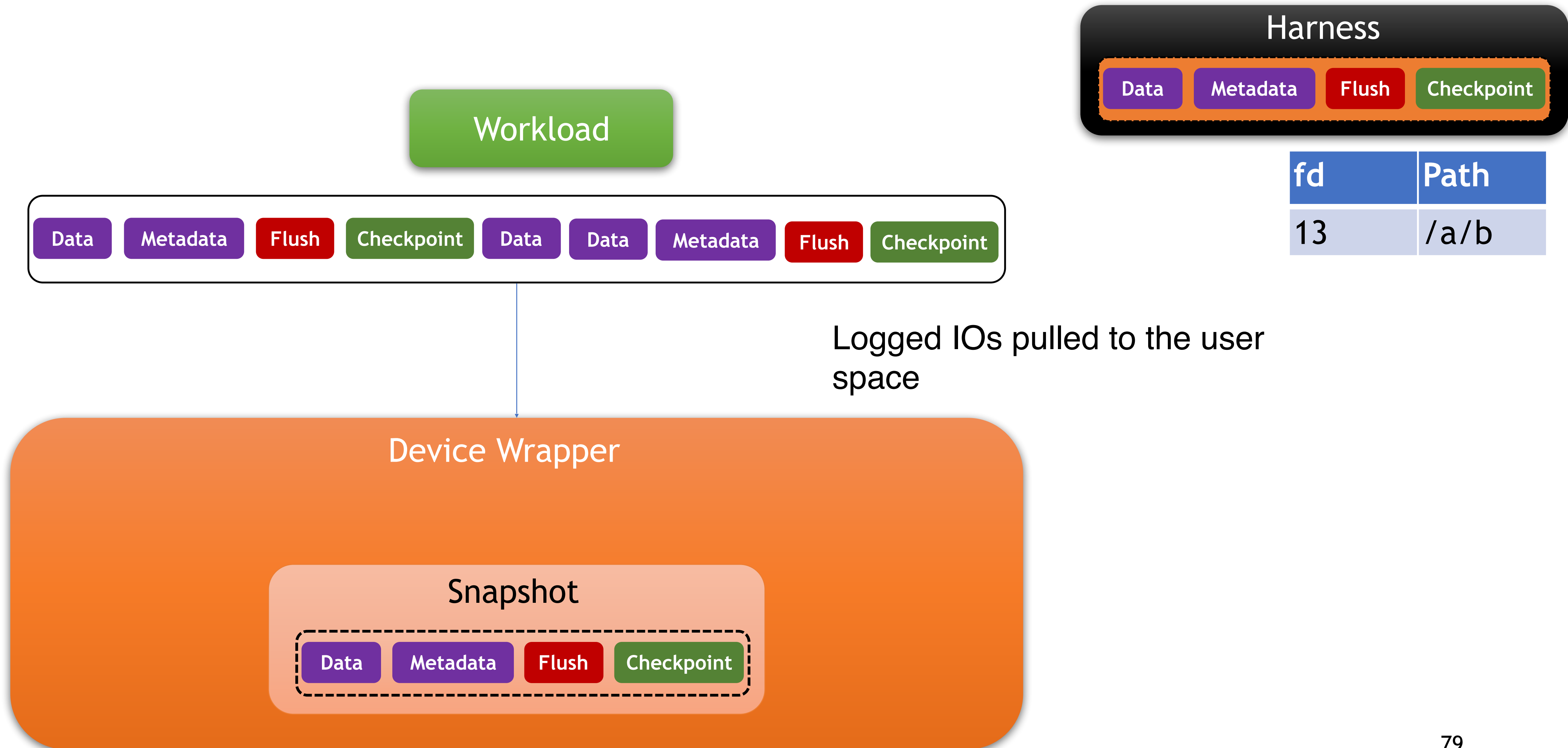
Harness

fd	Path
13	/a/b

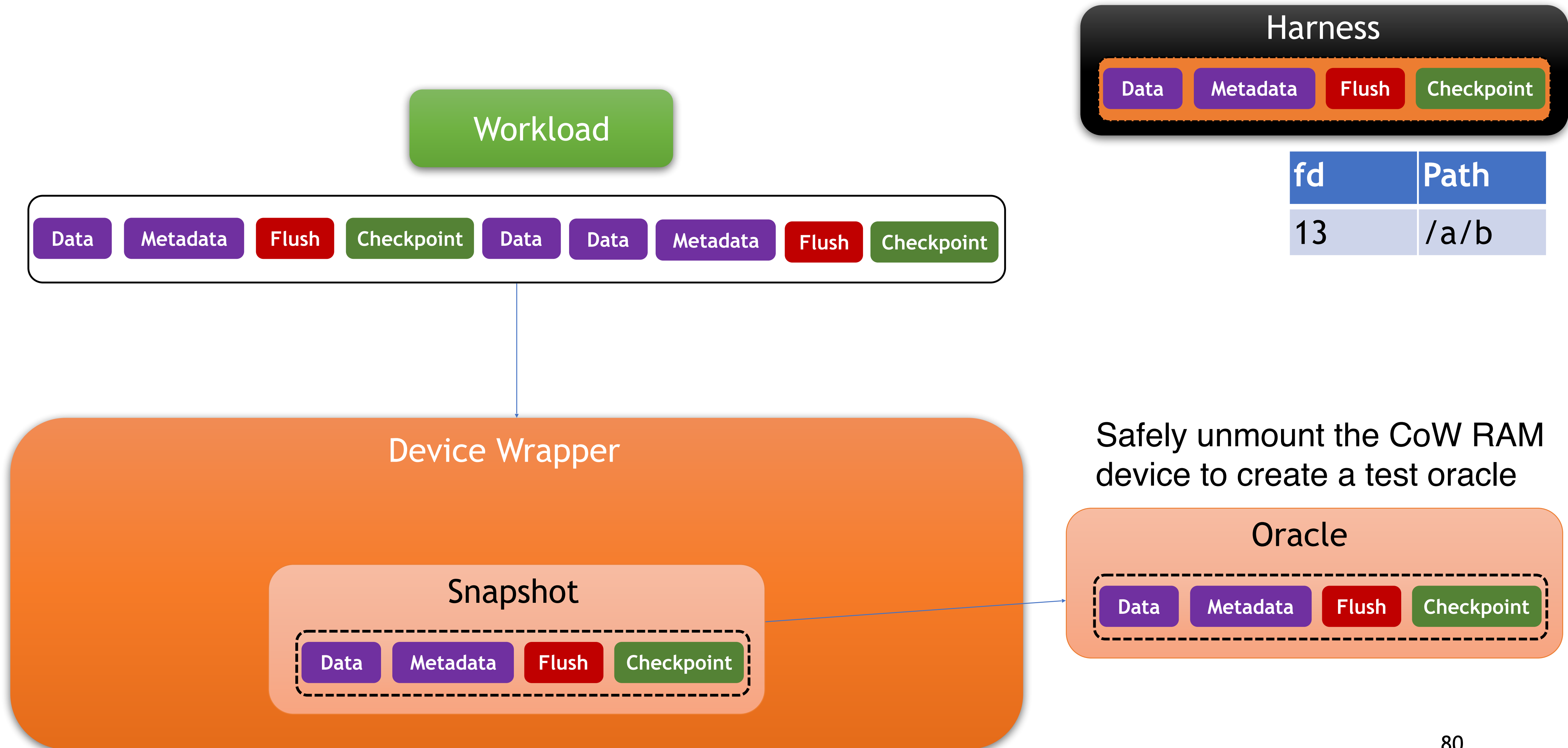
CrashMonkey in Action : Profiling



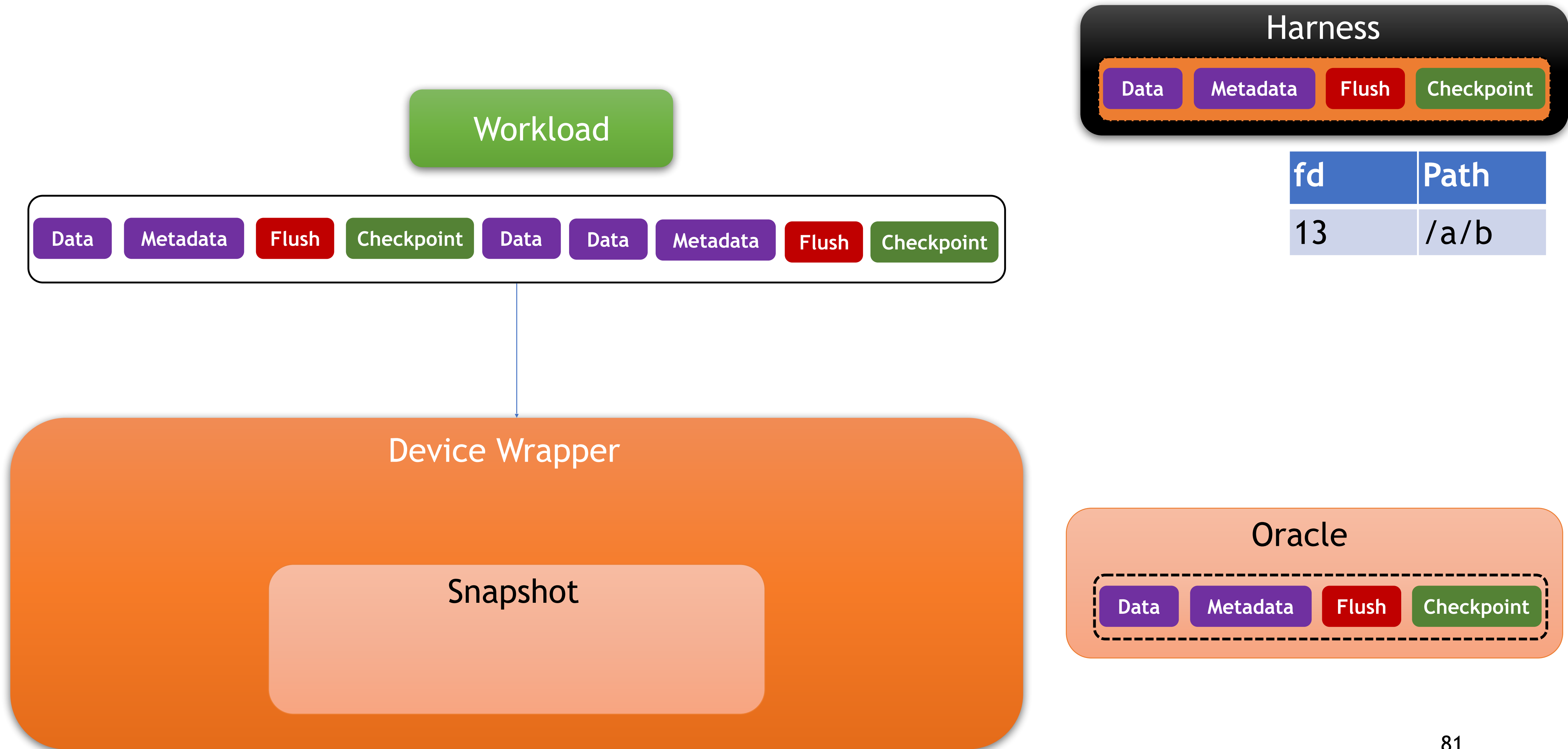
CrashMonkey in Action : Profiling



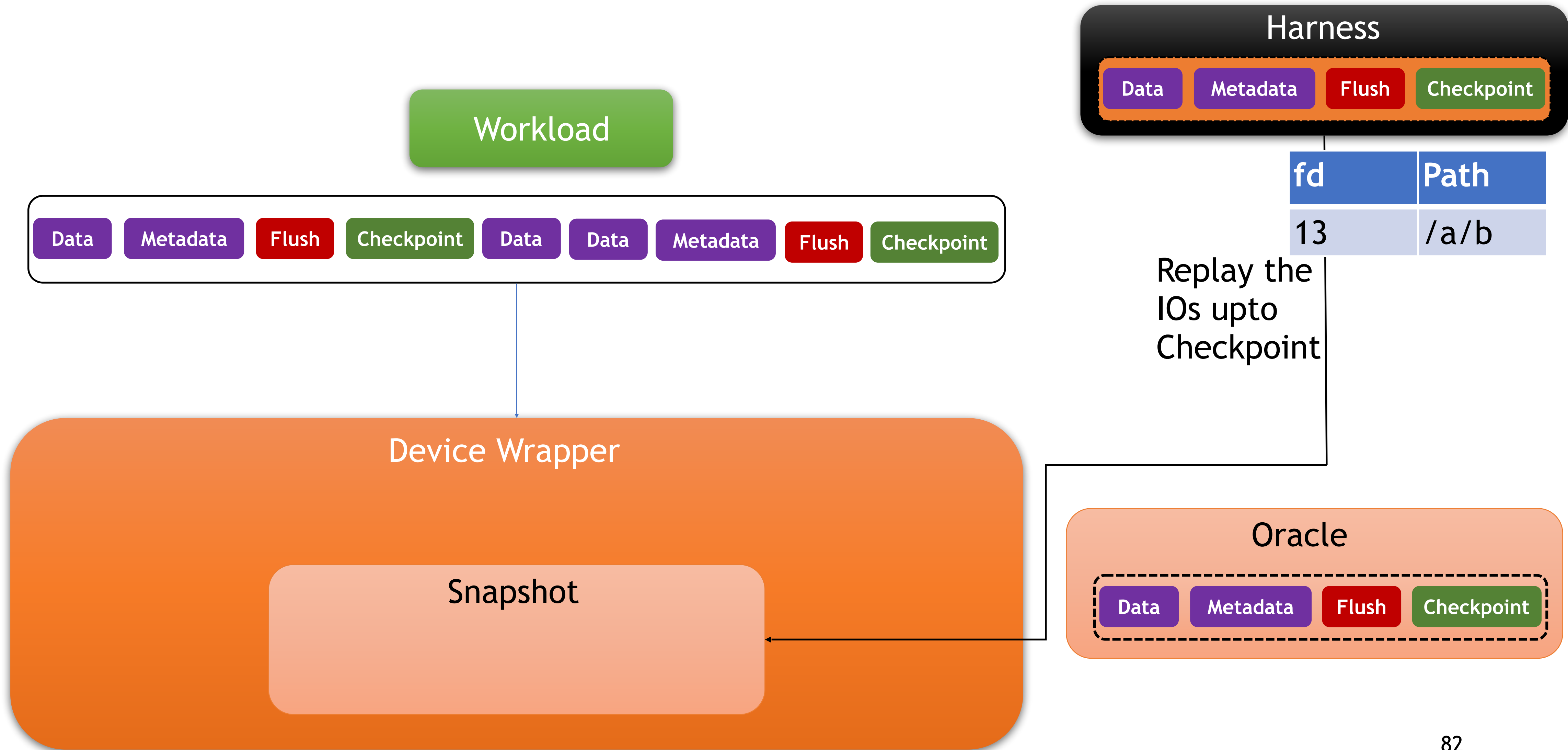
CrashMonkey in Action : Profiling



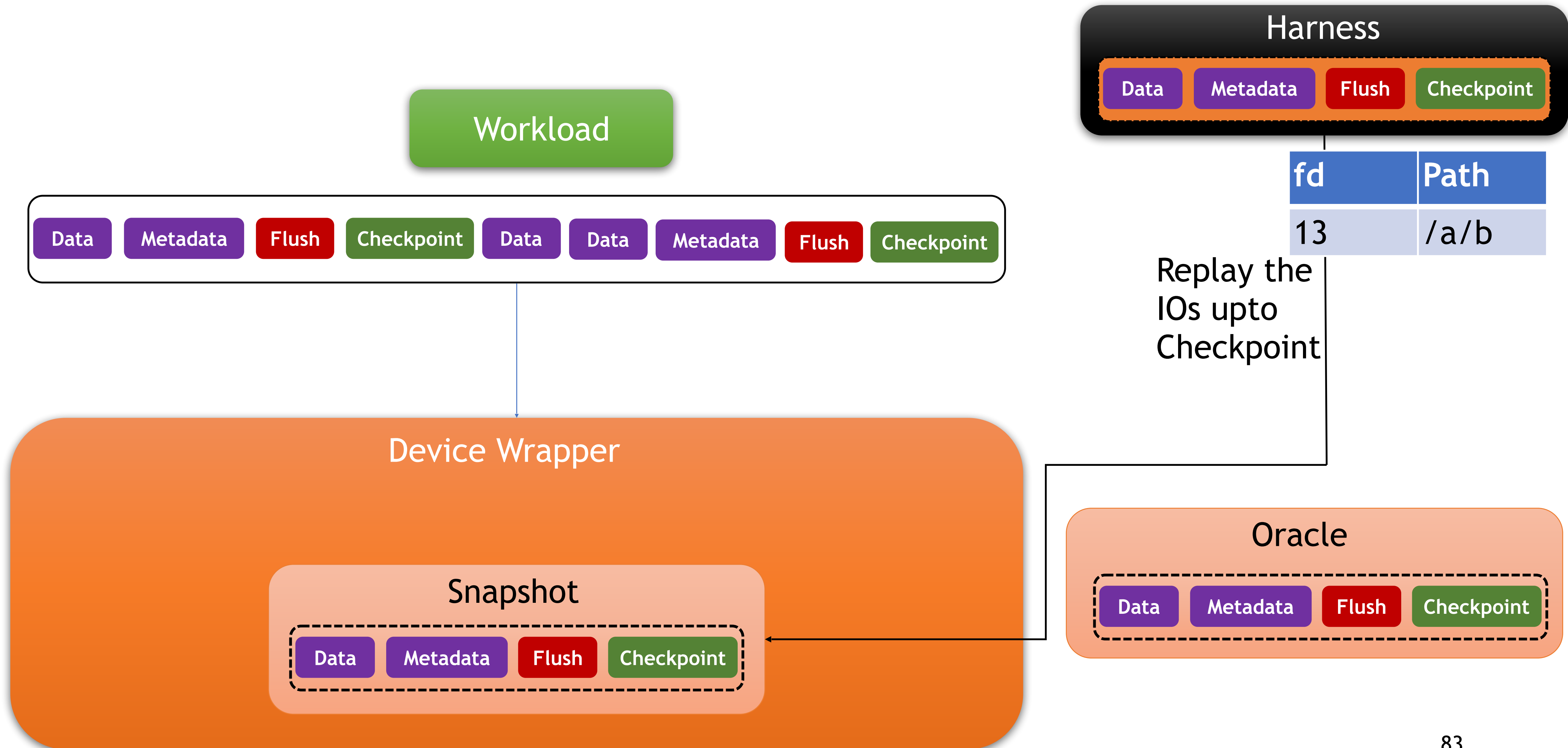
CrashMonkey in Action : Profiling



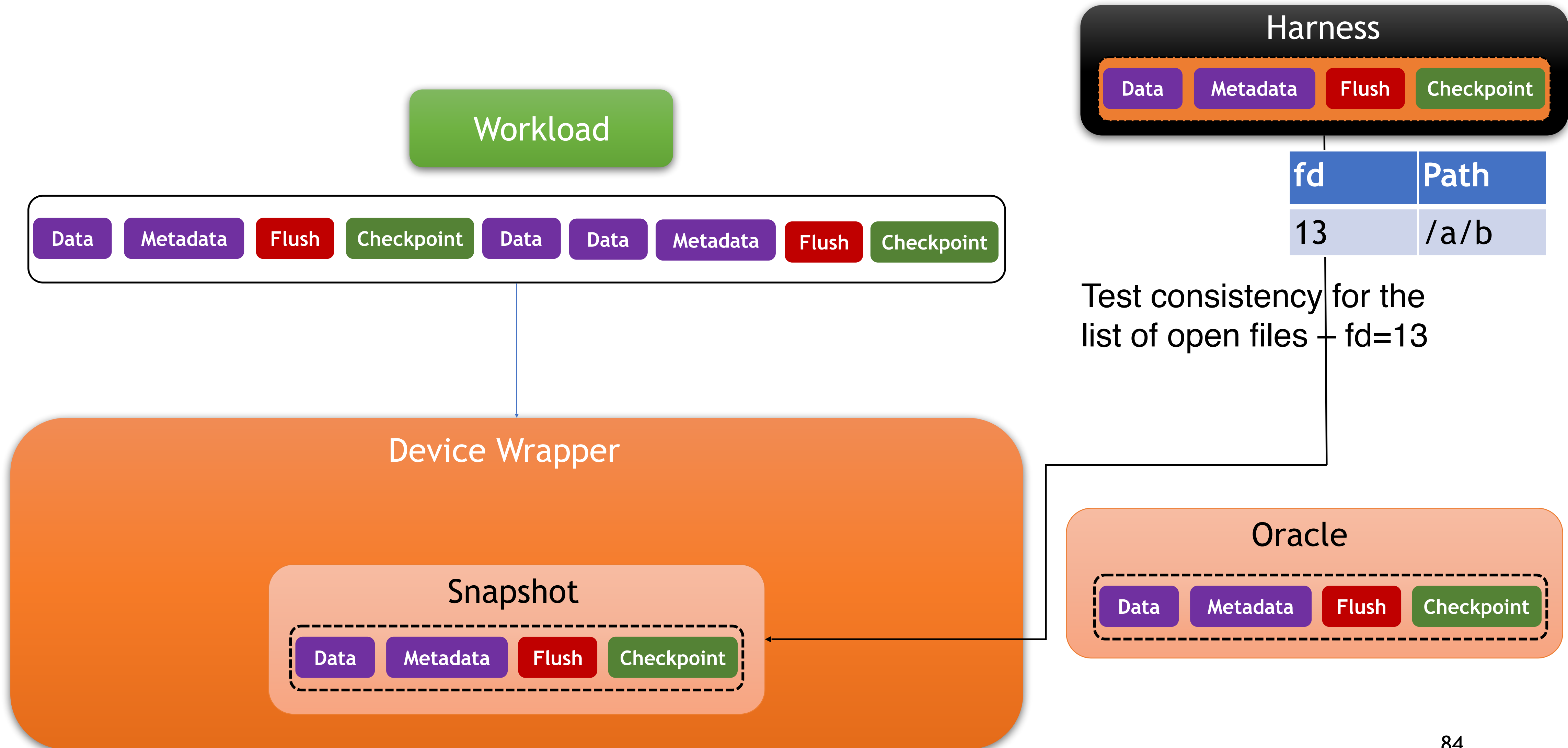
CrashMonkey in Action : Replay



CrashMonkey in Action : Replay



CrashMonkey in Action : Testing



Testing Strategy to find new bugs

- We test seq-1, seq-2 workloads on all filesystems : ext4, xfs, f2fs, btrfs
- We run all other workloads on btrfs and F2FS first.
 - For every workload that generated a bug, we run it on all other FS
- To run all workloads upto seq-3, you need to dedicate 2 days of compute per filesystem with (testing in parallel on 780 VM)

Results at a glance

Sequence Length	# workloads	# Bugs Reproduced	# Bugs found
Seq-1			
Seq-2			
Seq-3 metadata			
Seq-3 data			
Seq-3 nested			
Total			

- 25 million workloads
- Needs 15 days of testing on 780 VMs in parallel!

Results at a glance

Sequence Length	# workloads	# Bugs Reproduced	# Bugs found
Seq-1	300	3	3
Seq-2	254K	14	3
Seq-3 metadata	120K	5	2
Seq-3 data	1.5M	2	0
Seq-3 nested	1.5M	2	2
Total	3.37M	26	10