

CRASH MONKEY & ACE

Systematically Testing File-System Crash Consistency

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

[Published at OSDI 2018]



Crashes



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

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



Rename atomicity bug in btrfs

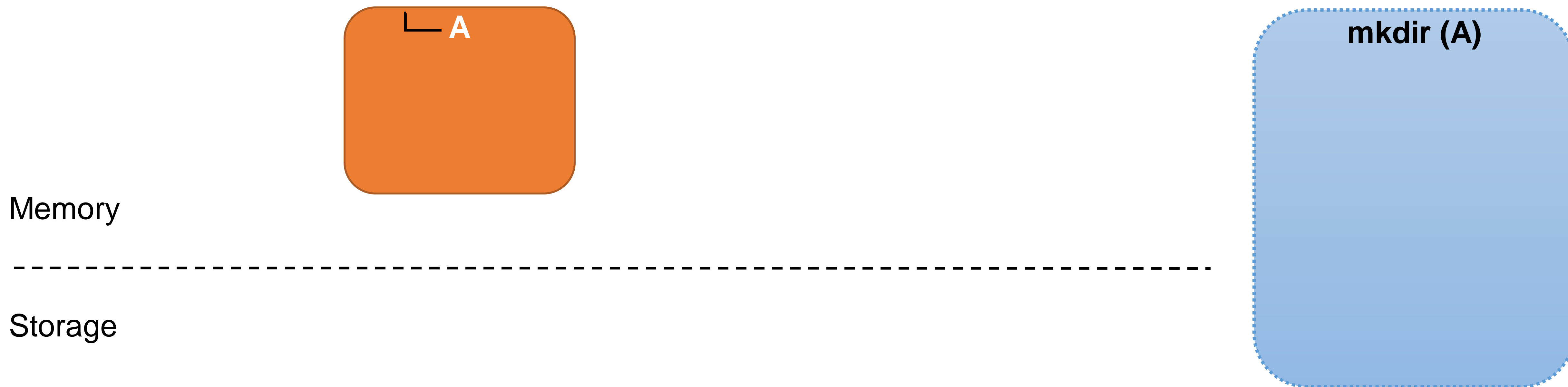
Memory



Storage



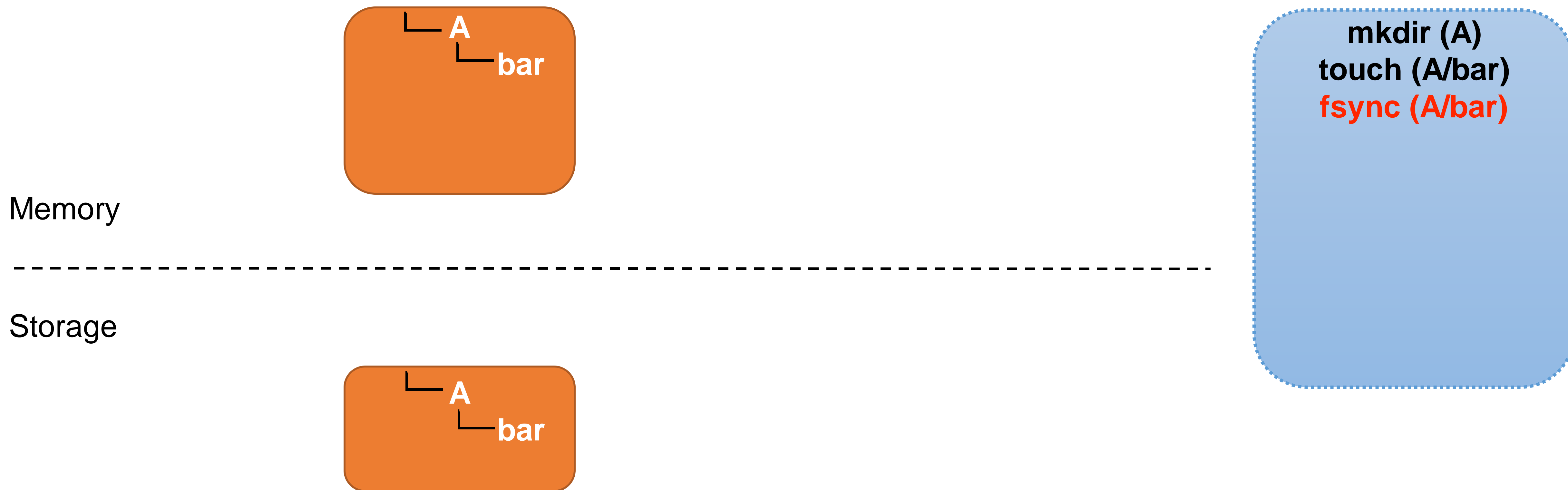
Rename atomicity bug in btrfs



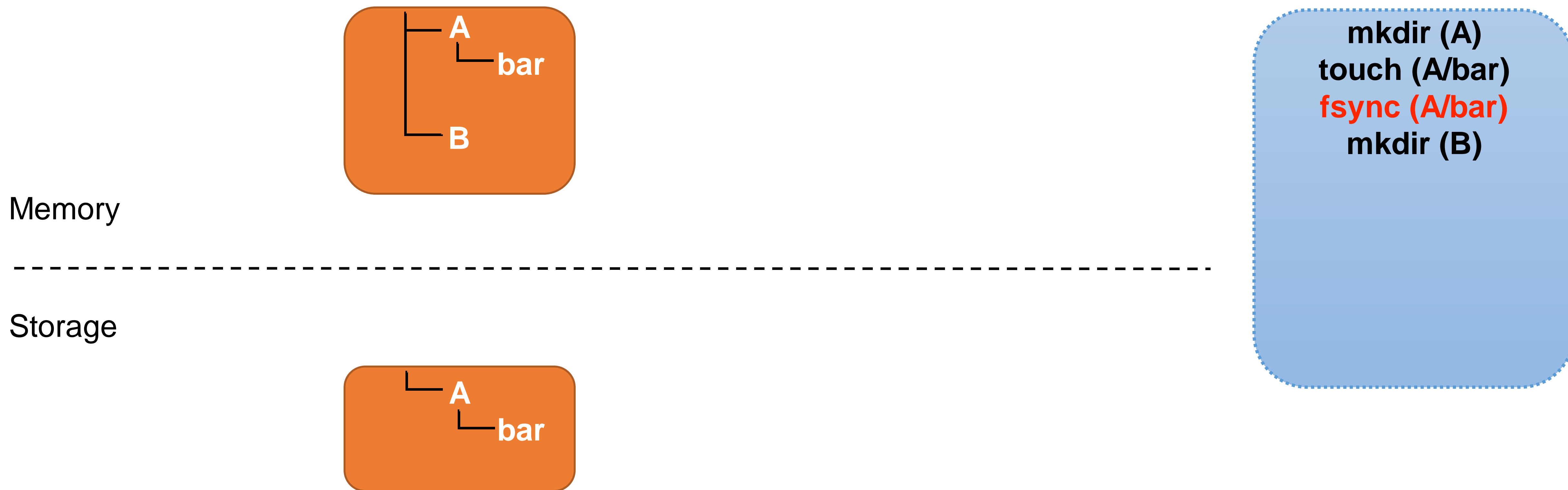
Rename atomicity bug in btrfs



Rename atomicity bug in btrfs

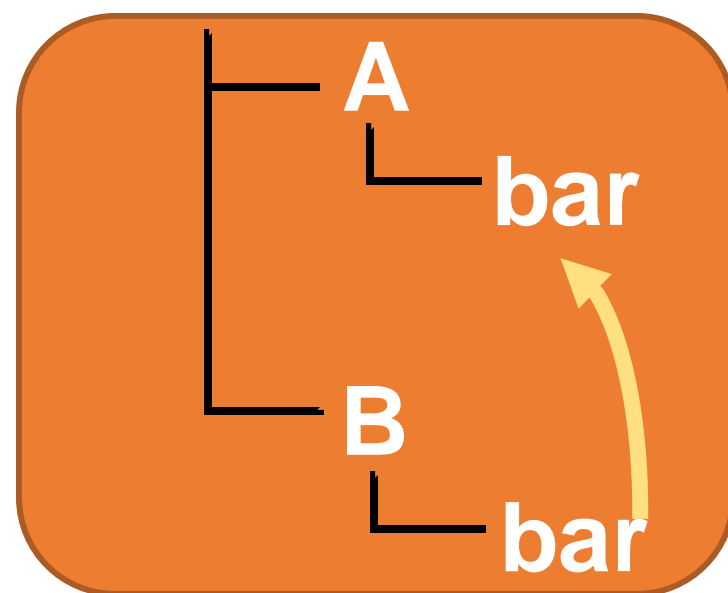


Rename atomicity bug in btrfs

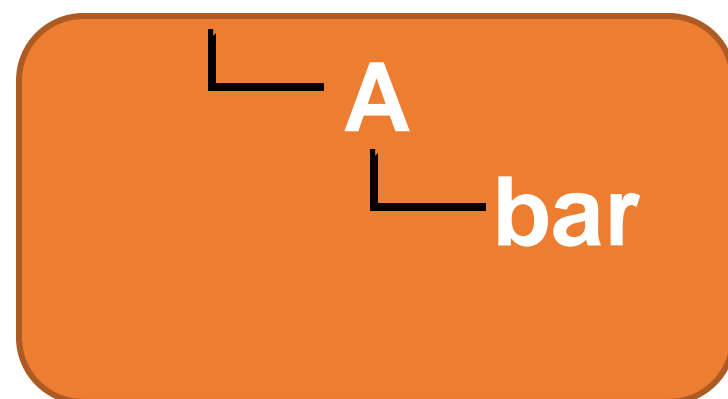


Rename atomicity bug in btrfs

Memory

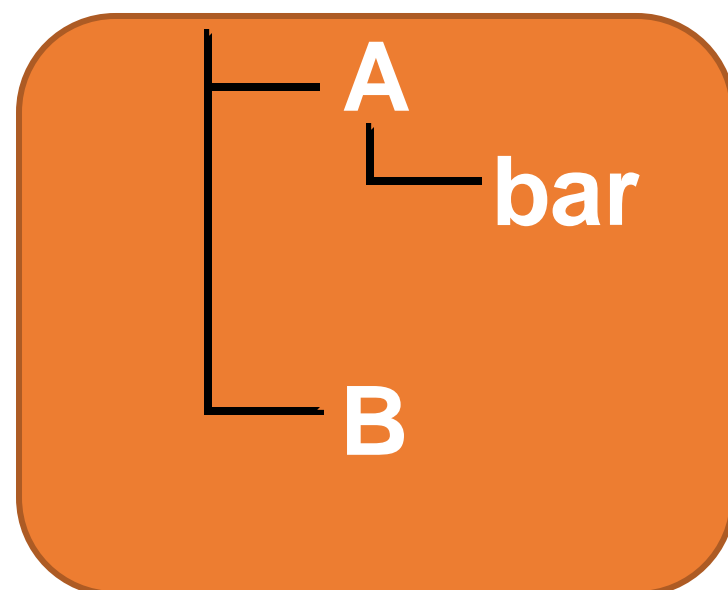


Storage

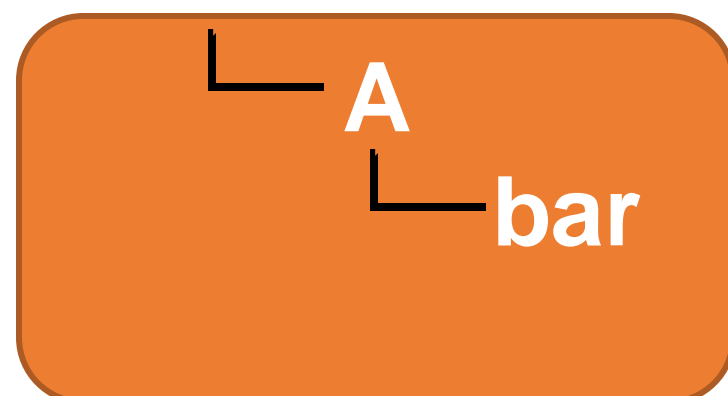


Rename atomicity bug in btrfs

Memory



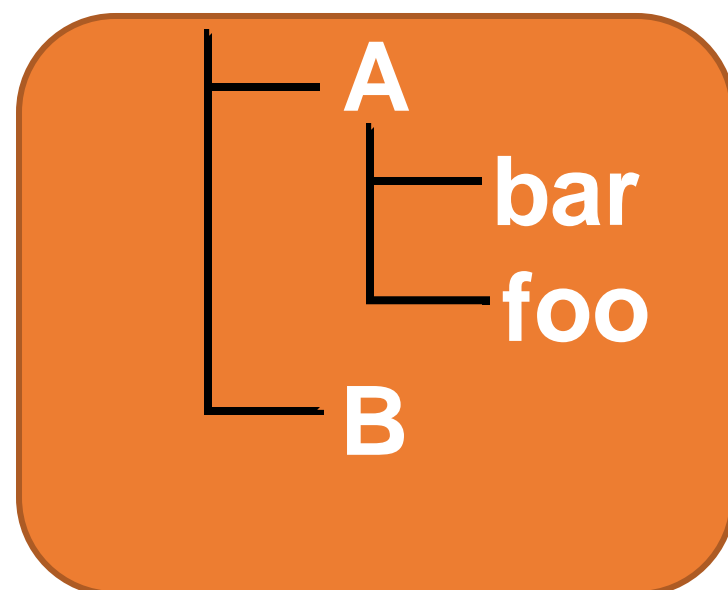
Storage



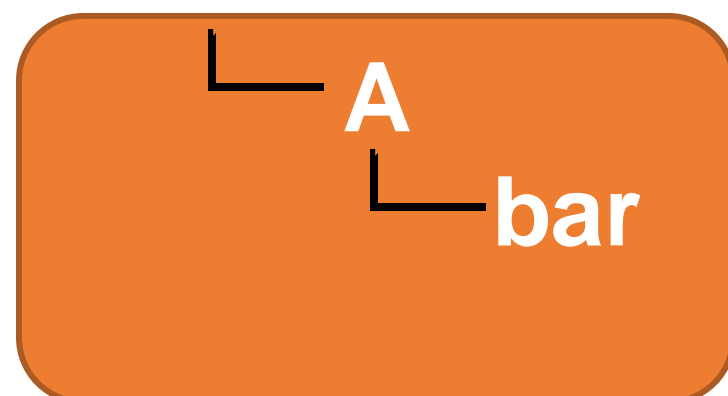
mkdir (A)
touch (A/bar)
fsync (A/bar)
mkdir (B)
touch (B/bar)
rename (B/bar, A/bar)

Rename atomicity bug in btrfs

Memory



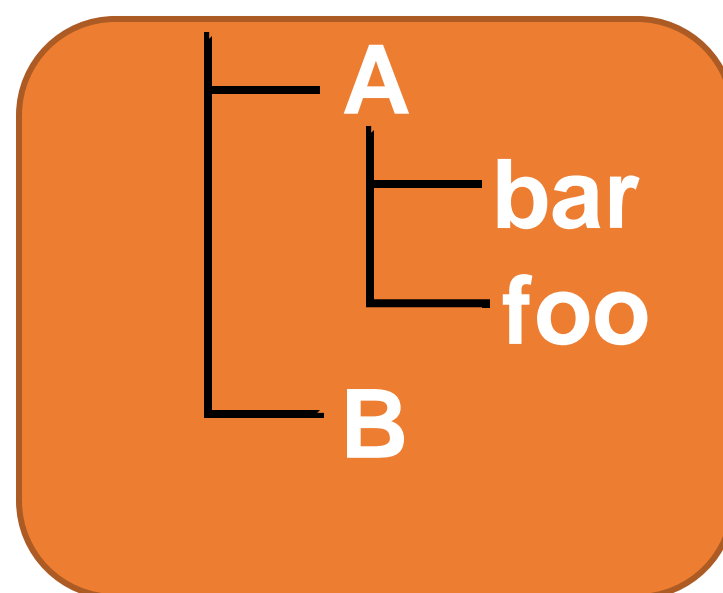
Storage



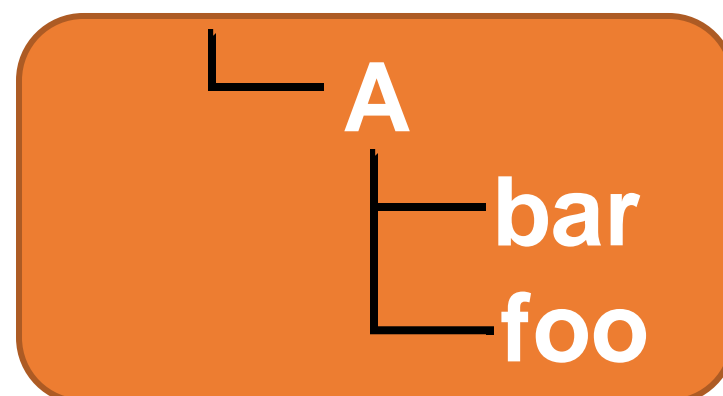
mkdir (A)
touch (A/bar)
fsync (A/bar)
mkdir (B)
touch (B/bar)
rename (B/bar, A/bar)
touch (A/foo)

Rename atomicity bug in btrfs

Memory

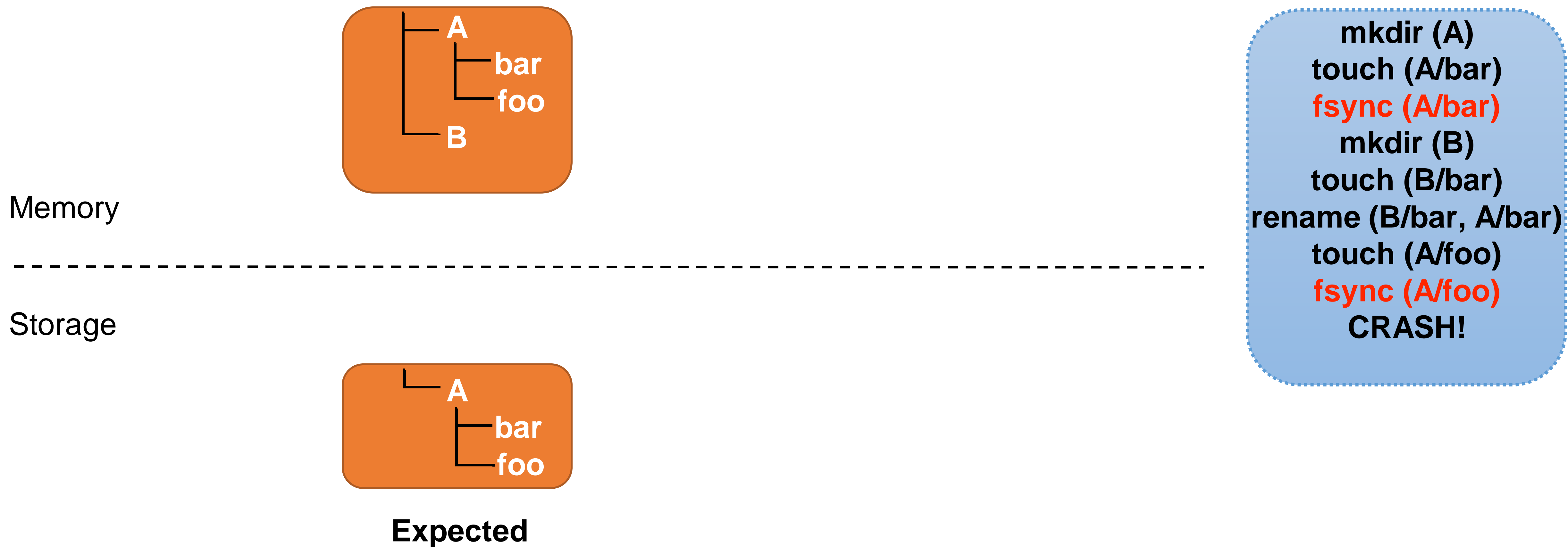


Storage

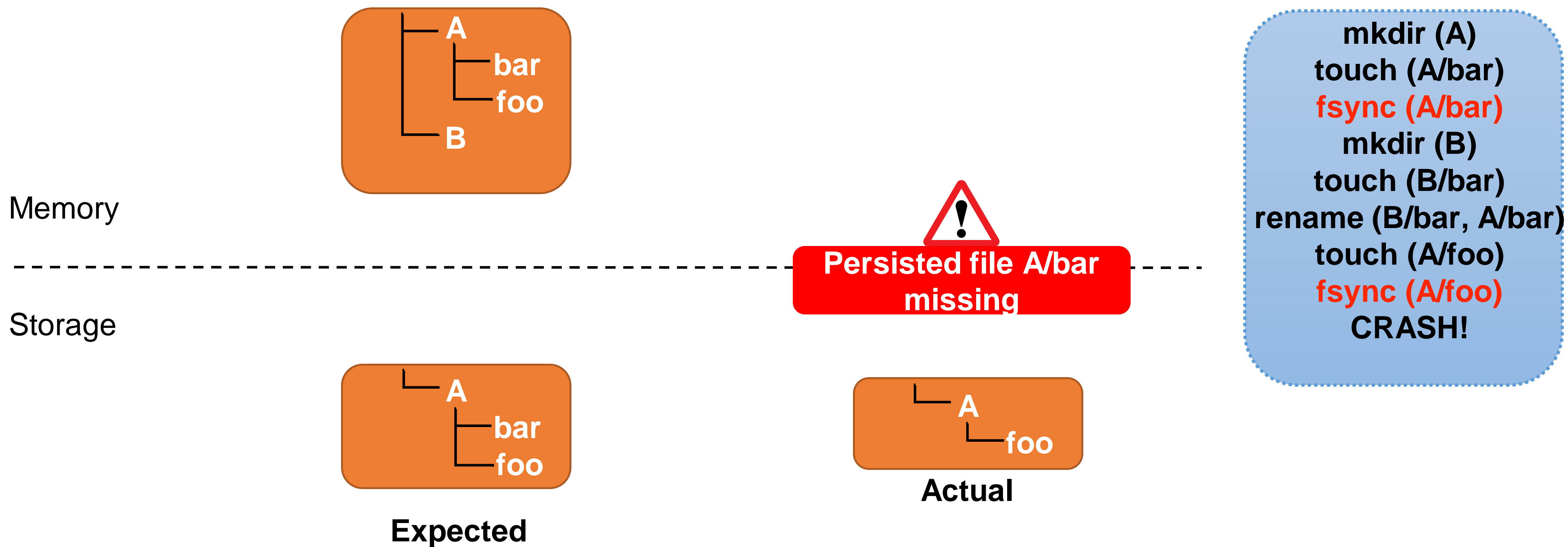


mkdir (A)
touch (A/bar)
fsync (A/bar)
mkdir (B)
touch (B/bar)
rename (B/bar, A/bar)
touch (A/foo)
fsync (A/foo)

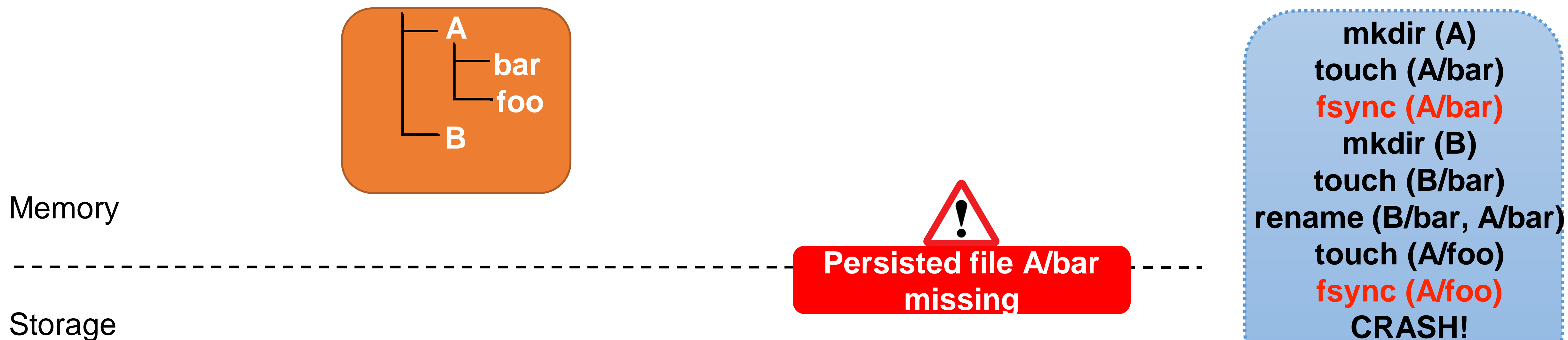
Rename atomicity bug in btrfs



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**

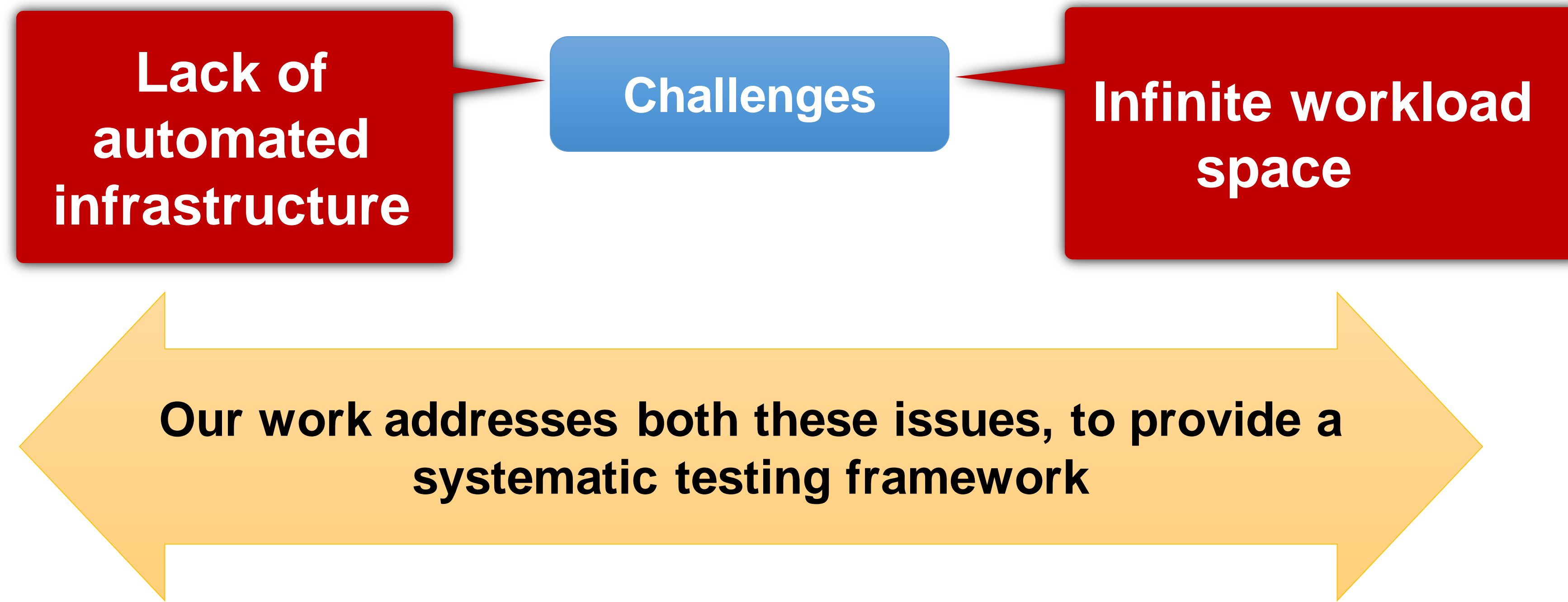
Annotate filesystems
Hard to do for existing FS

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

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

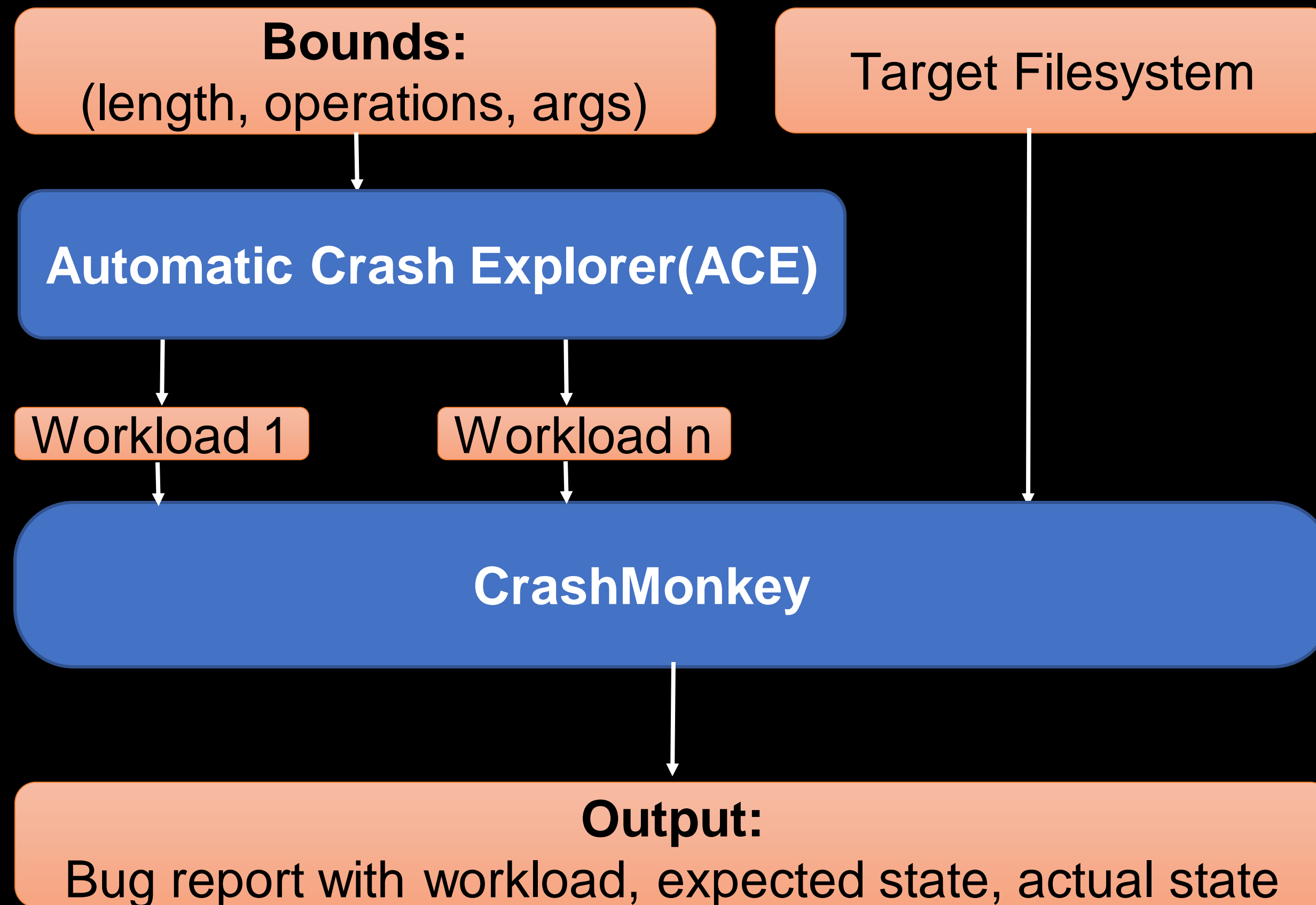
Challenges with systematic testing

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)

CrashMonkey and Ace : Features

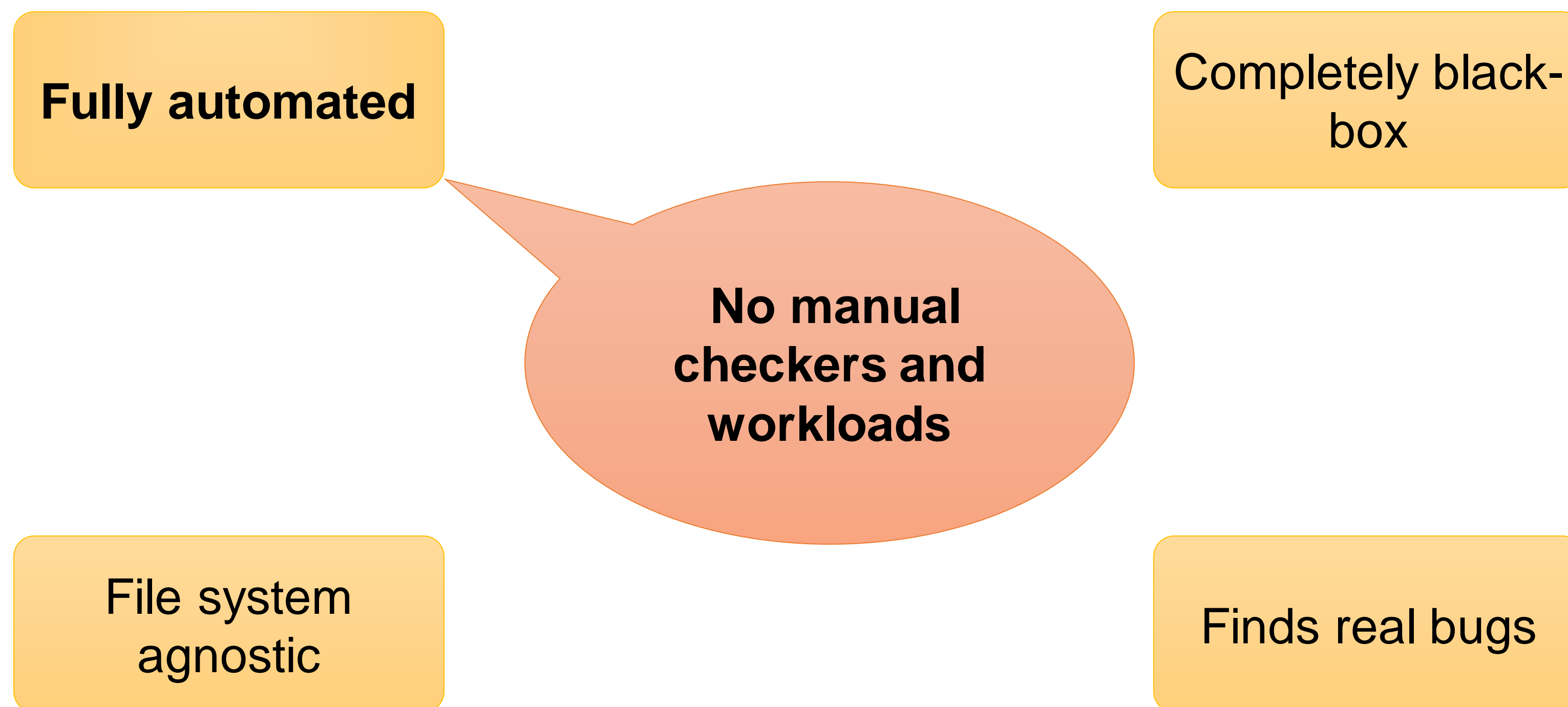
Fully automated

**Completely
black-box**

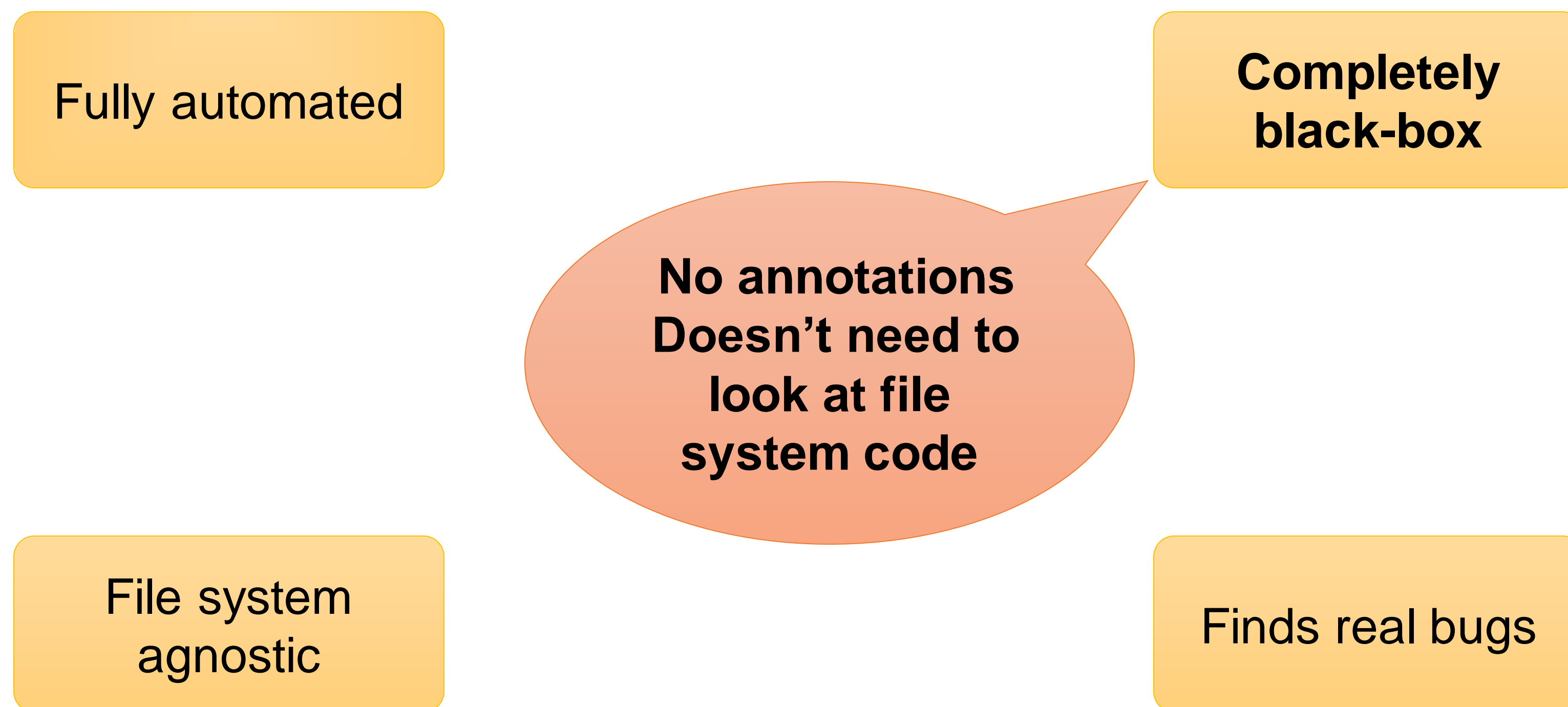
**File system
agnostic**

Finds real bugs

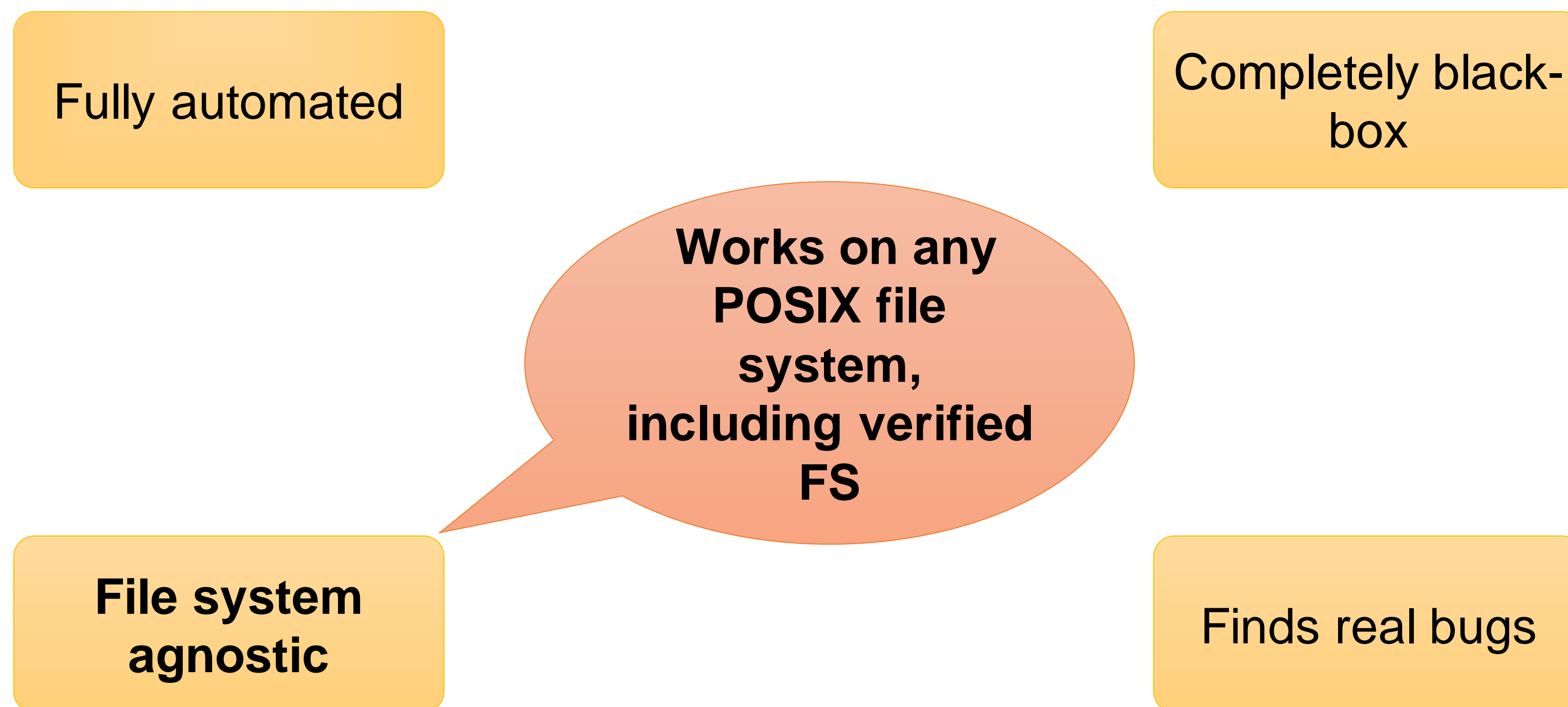
CrashMonkey and Ace : Features



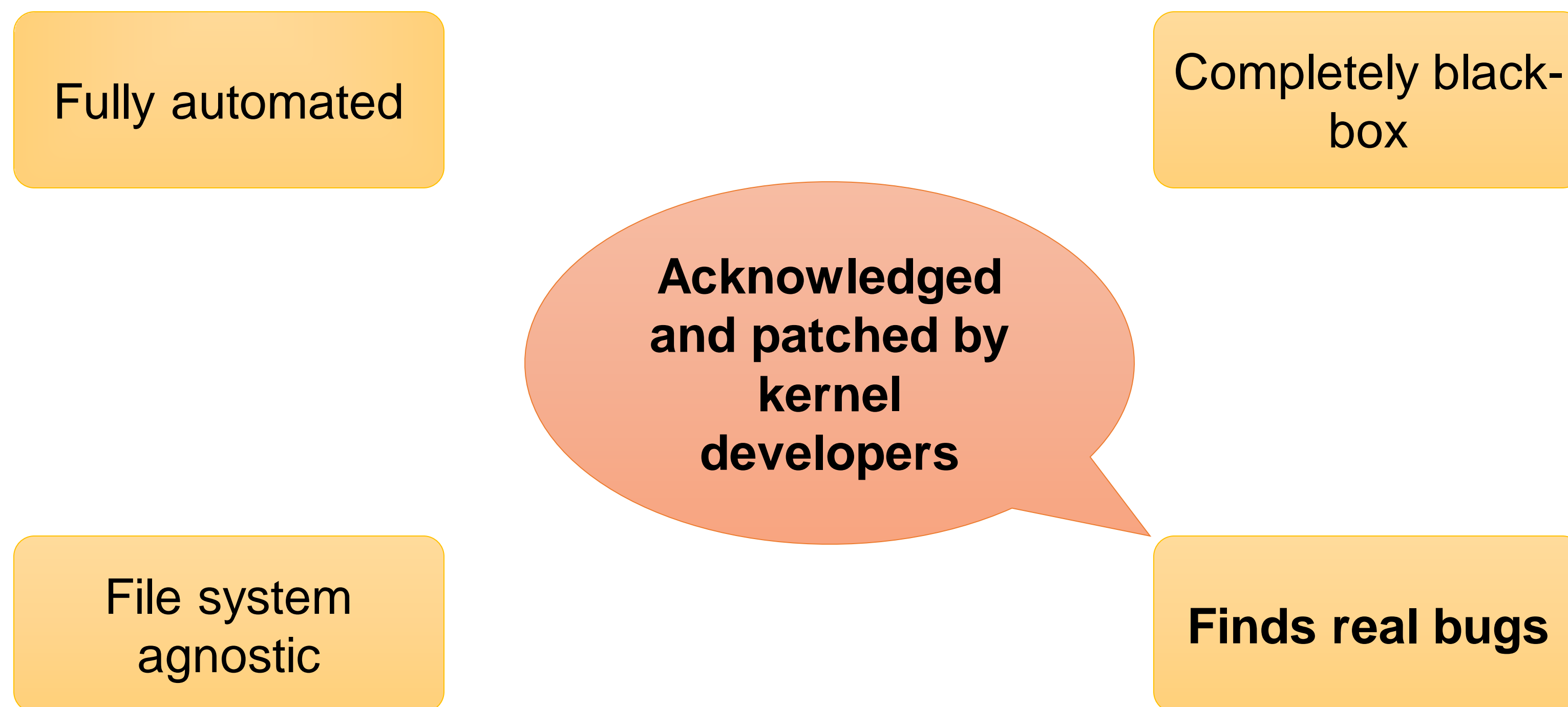
CrashMonkey and Ace : Features



CrashMonkey and Ace : Features



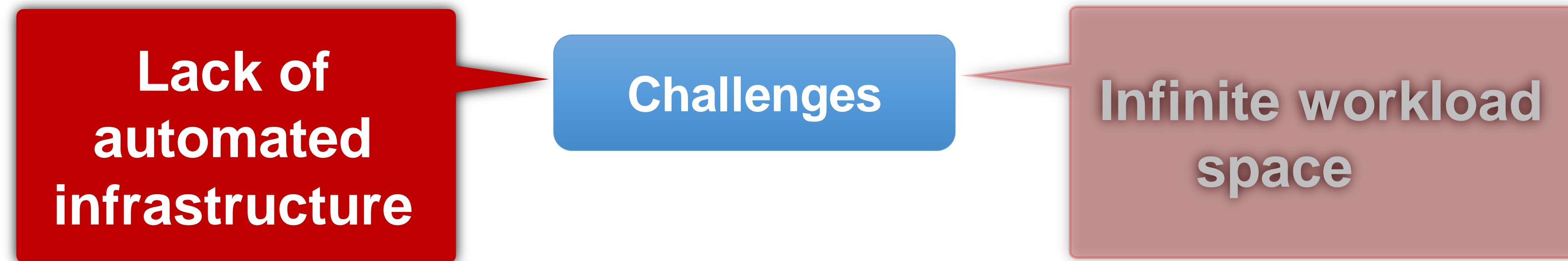
CrashMonkey and Ace : Features



Outline

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

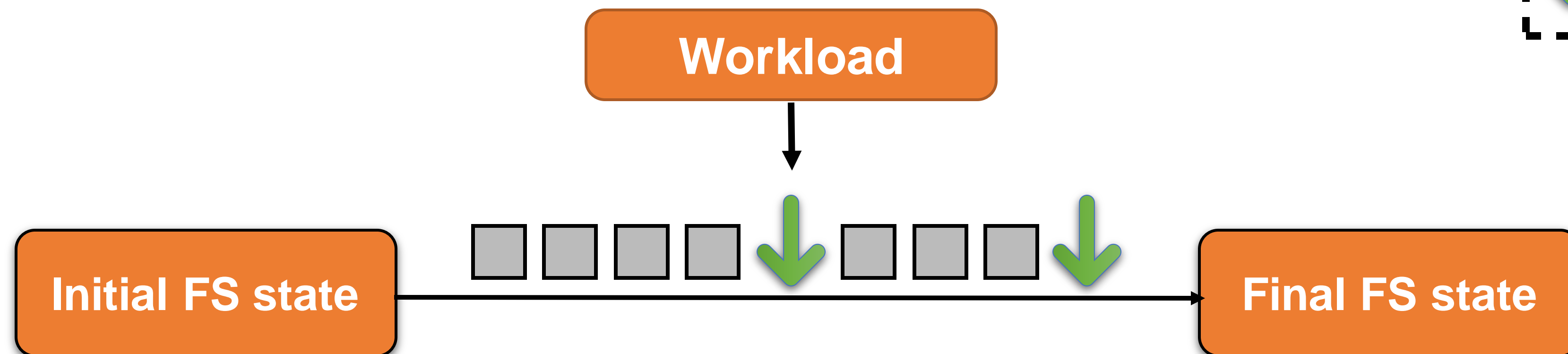
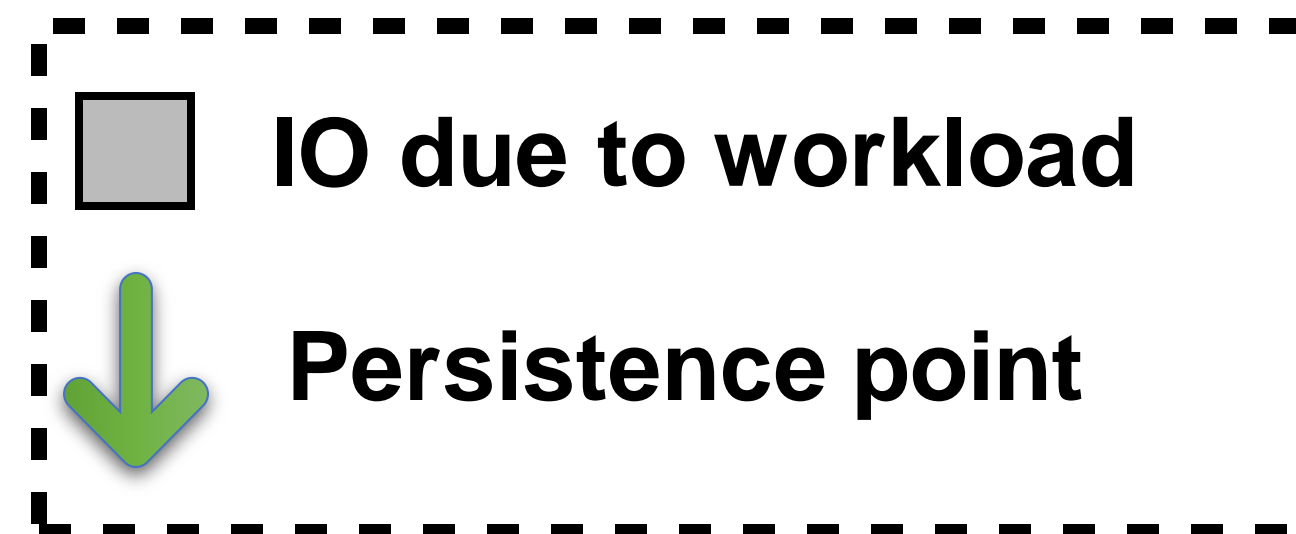
Challenges with systematic testing



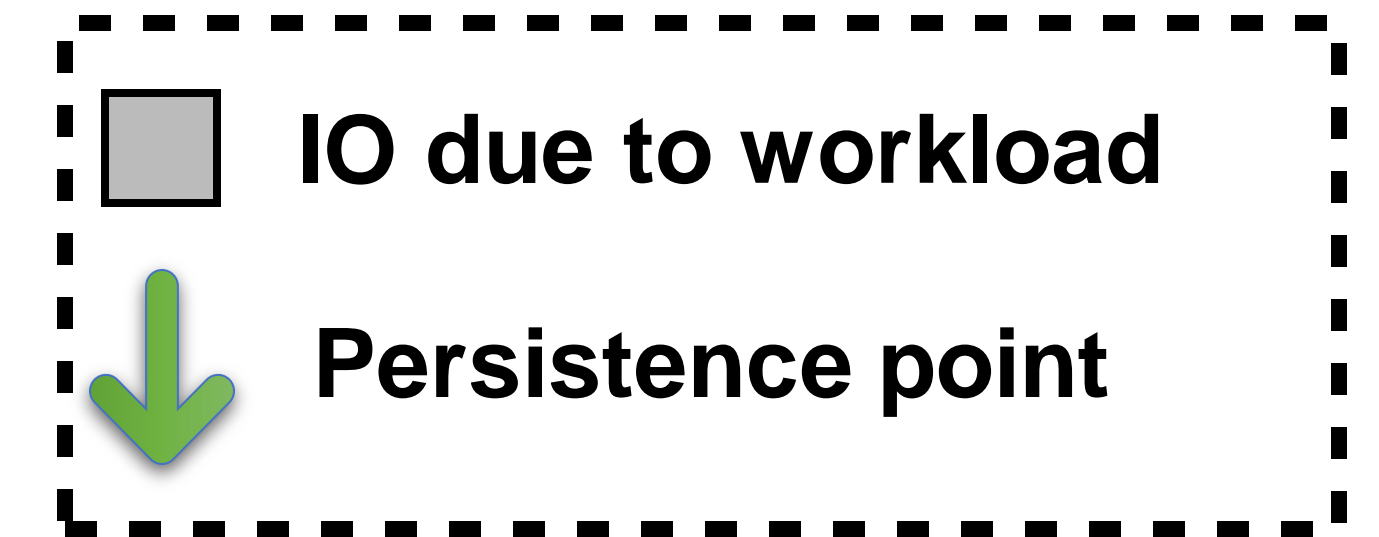
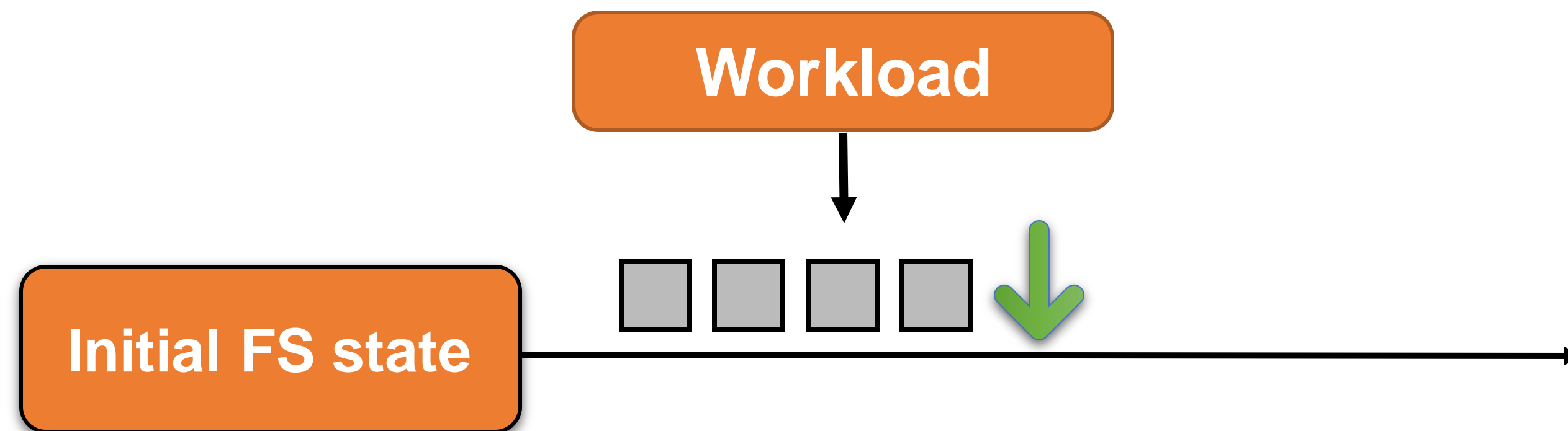
CrashMonkey

- 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

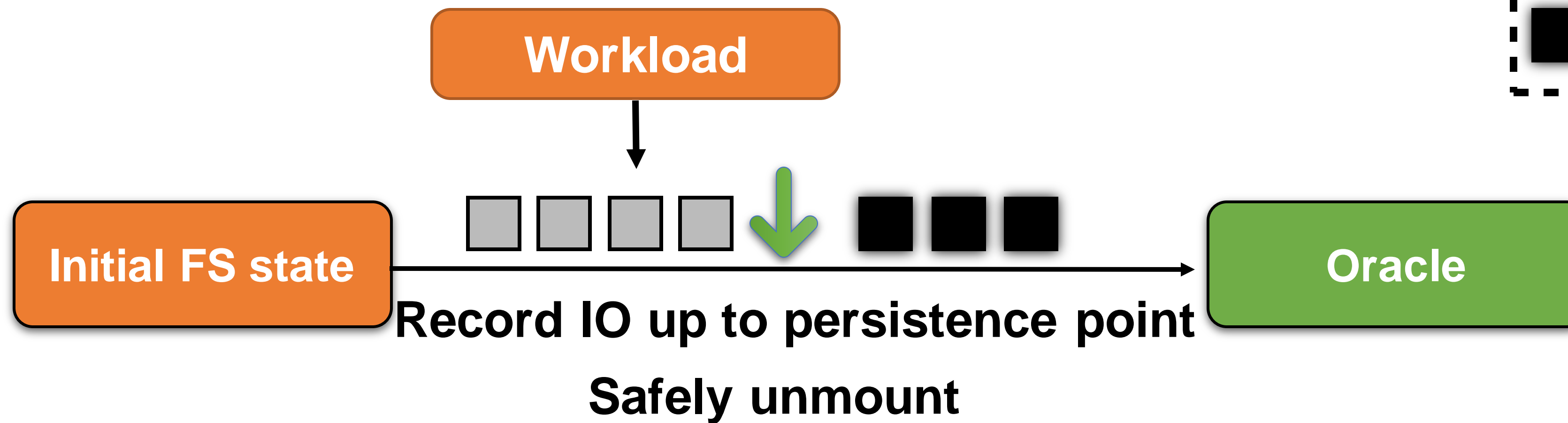
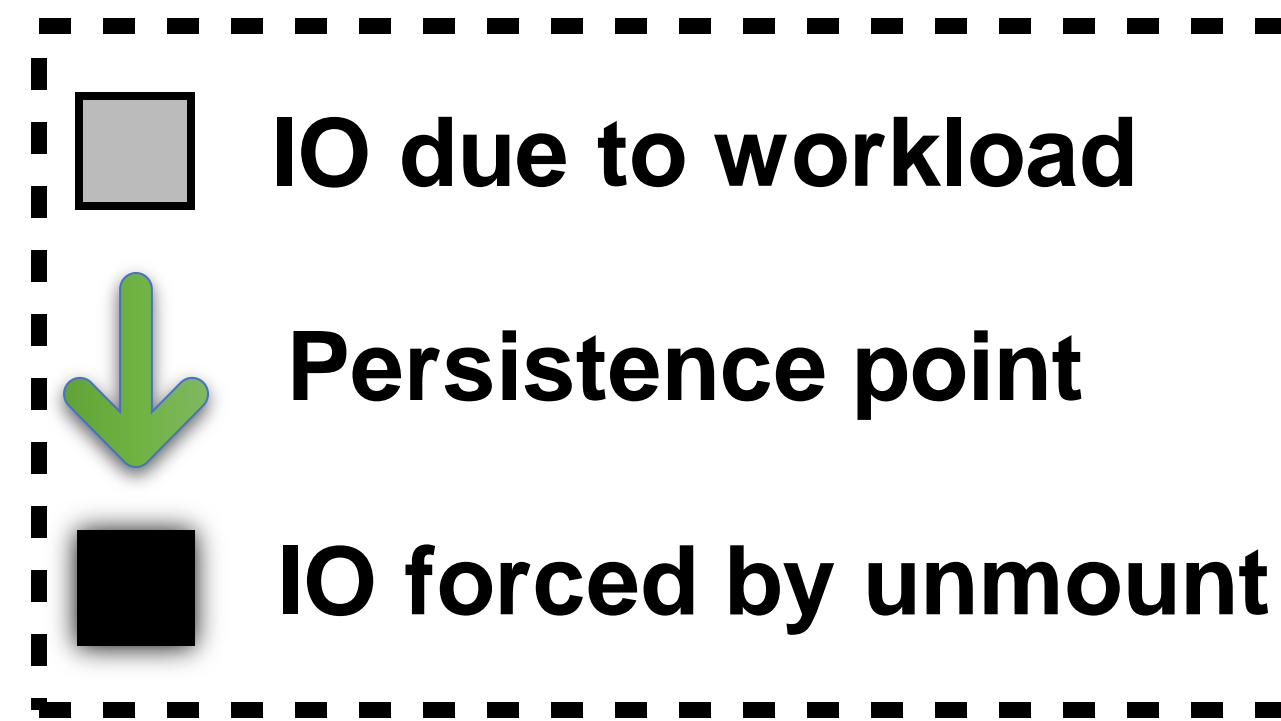
CrashMonkey in Action



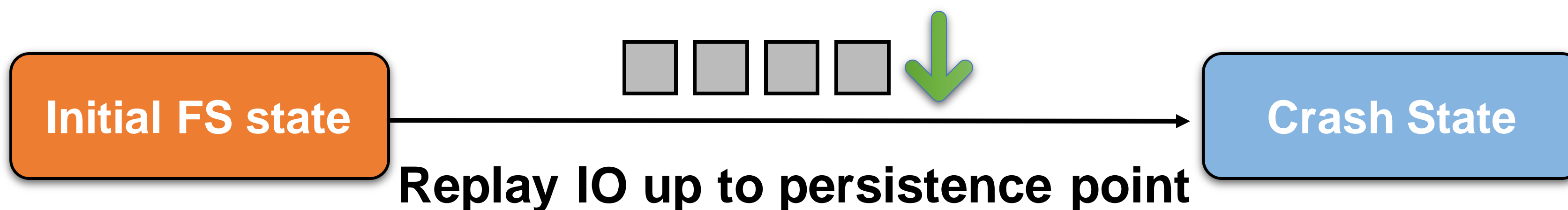
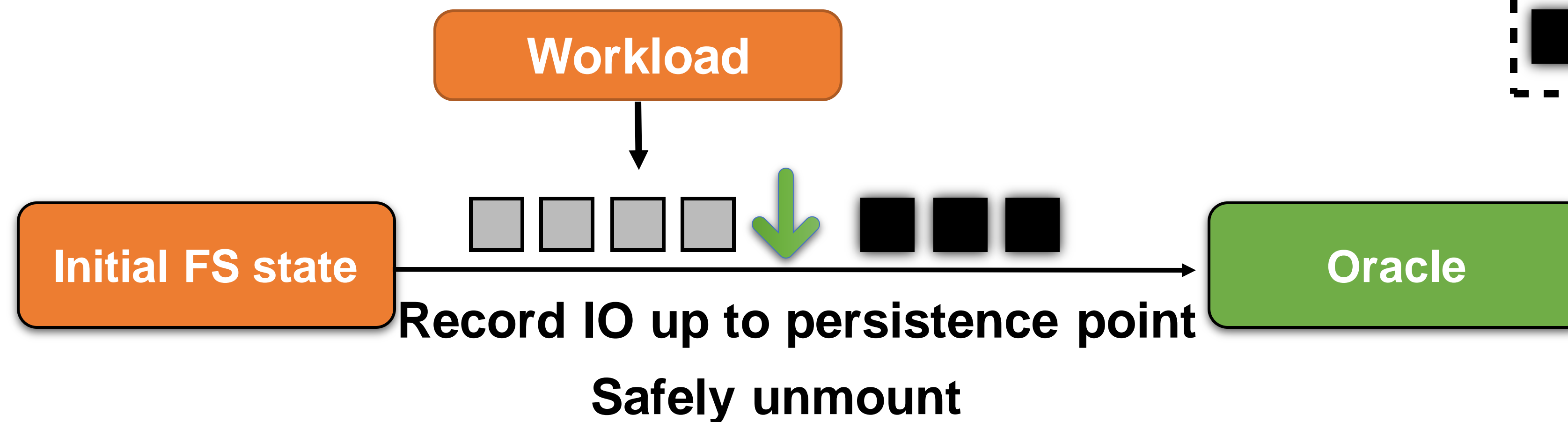
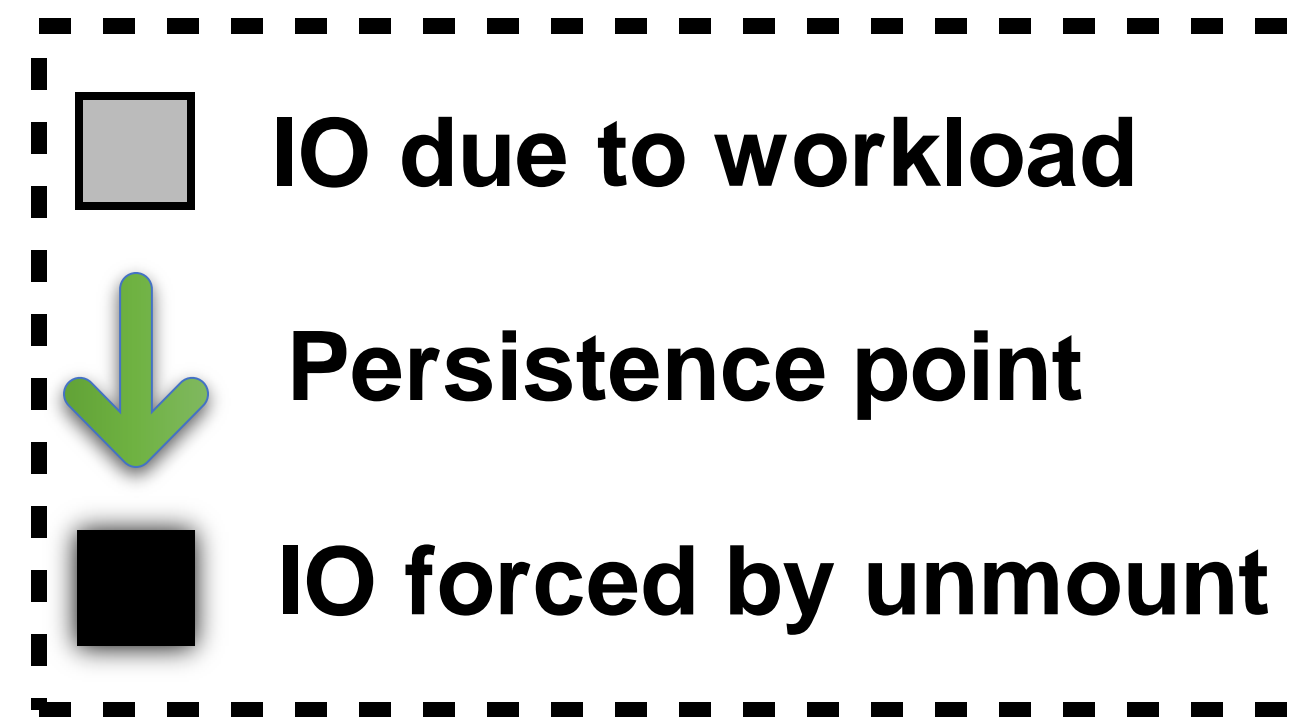
CrashMonkey in Action



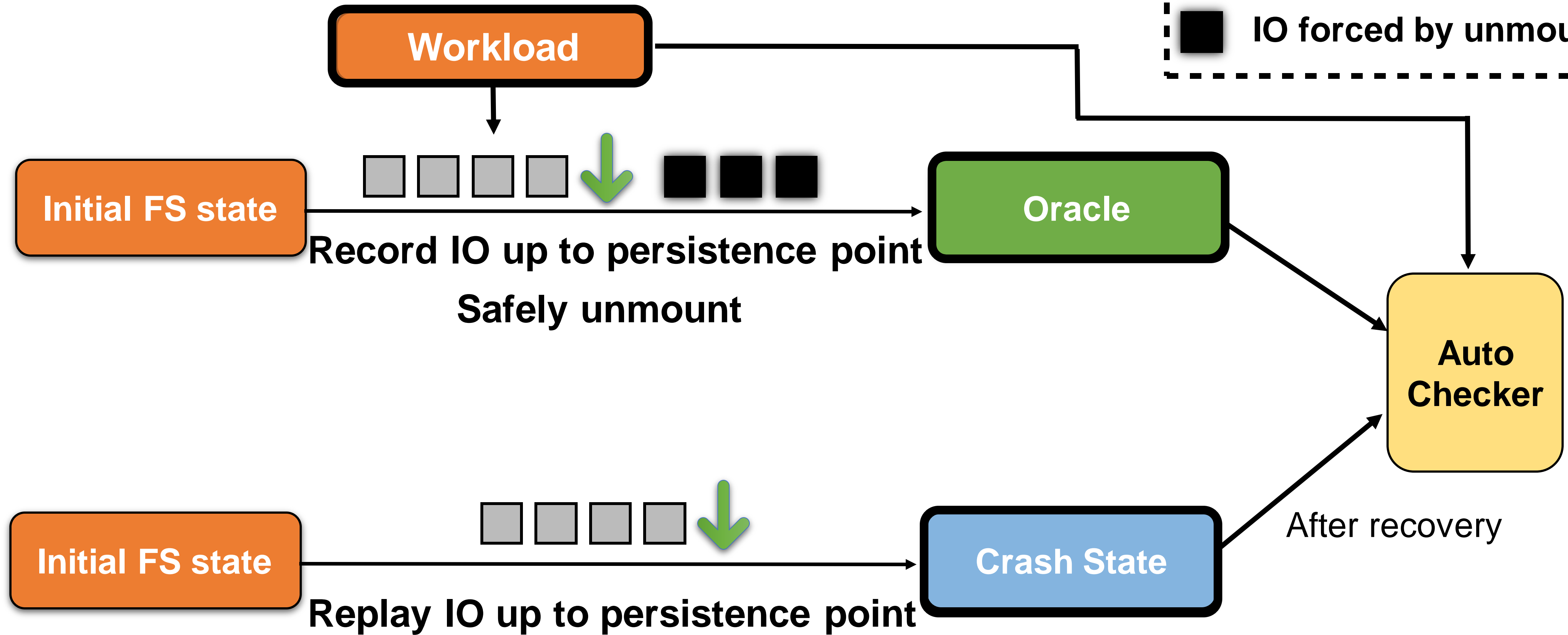
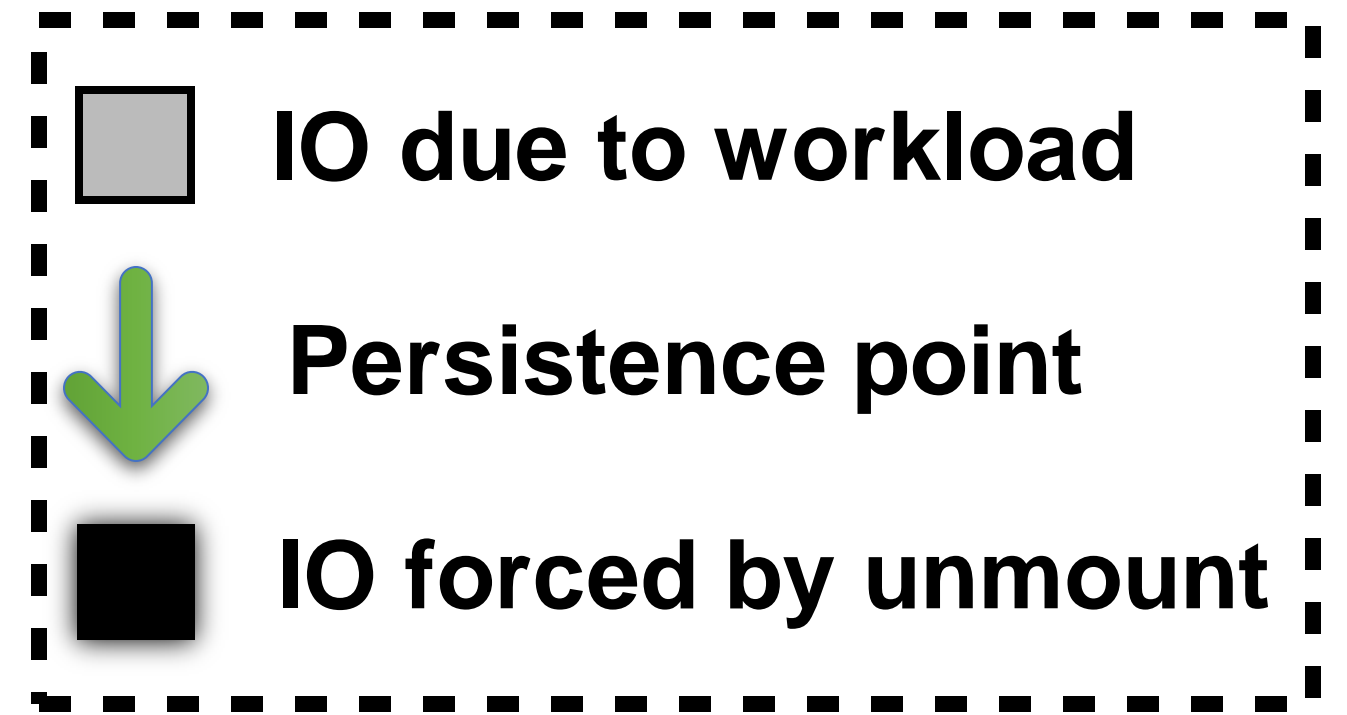
Phase 1 : Record IO



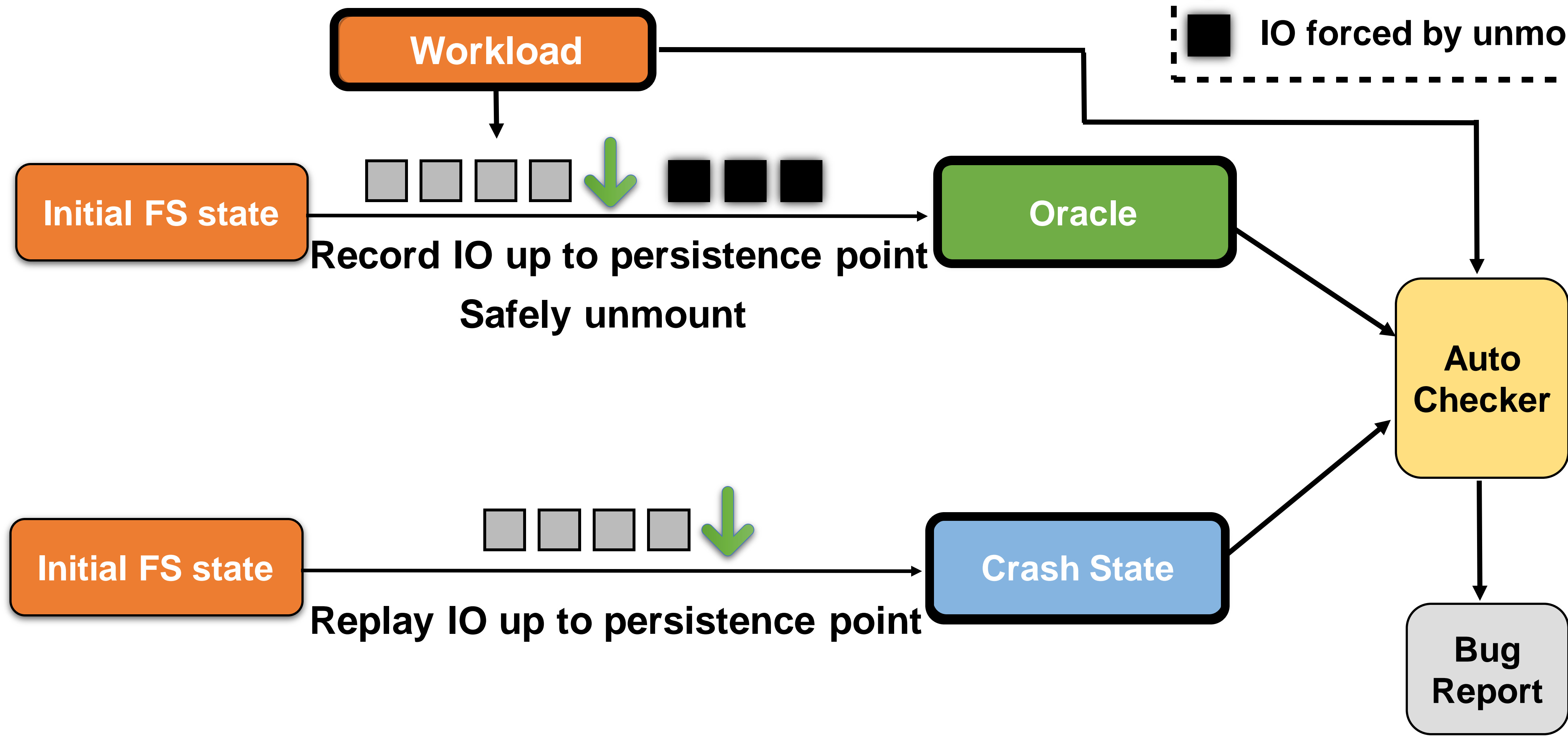
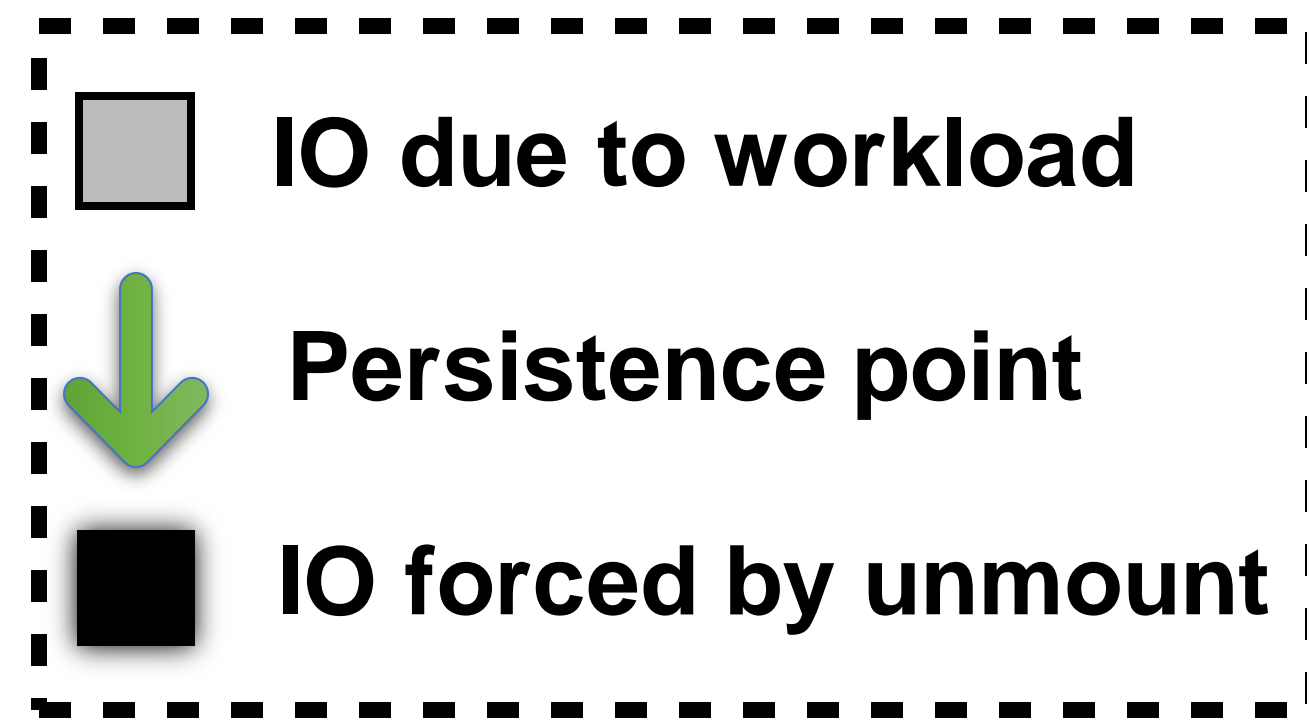
Phase 2 : Replay IO



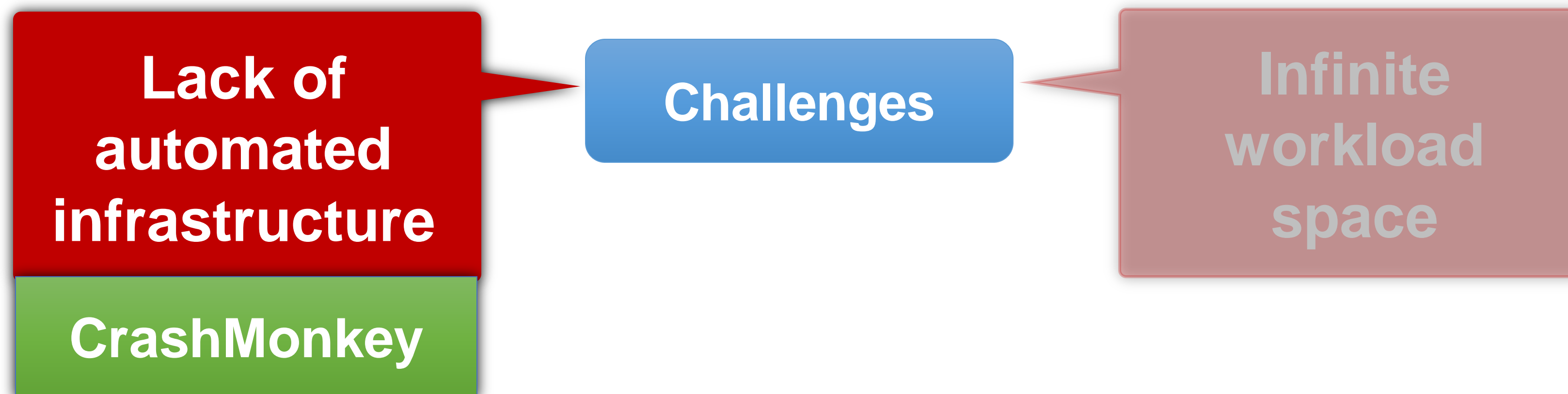
Phase 3 : Test for consistency



Phase 3 : Test for consistency



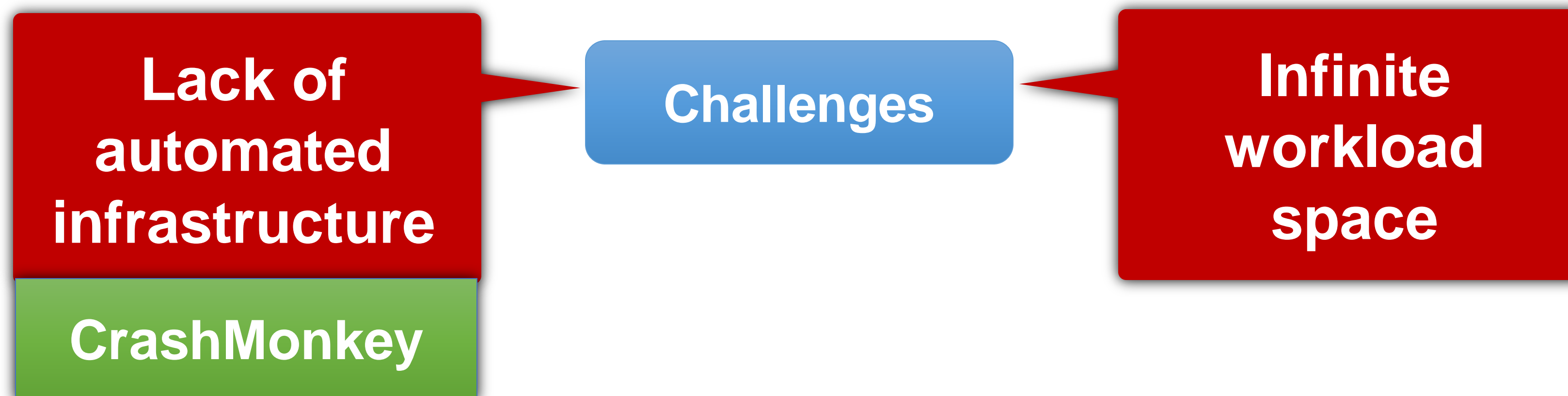
Challenges with Systematic Testing



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

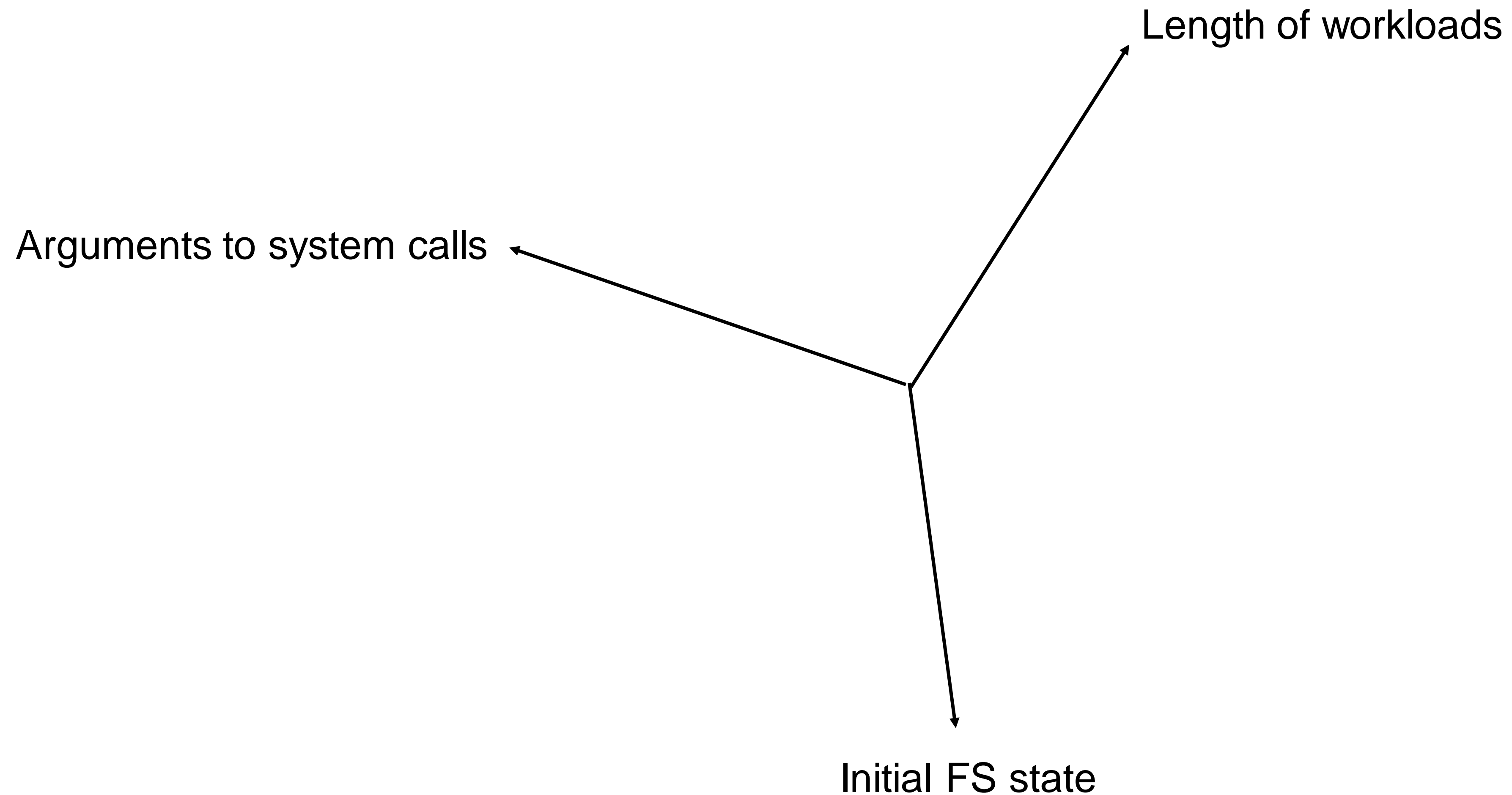
Challenges:

- 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?

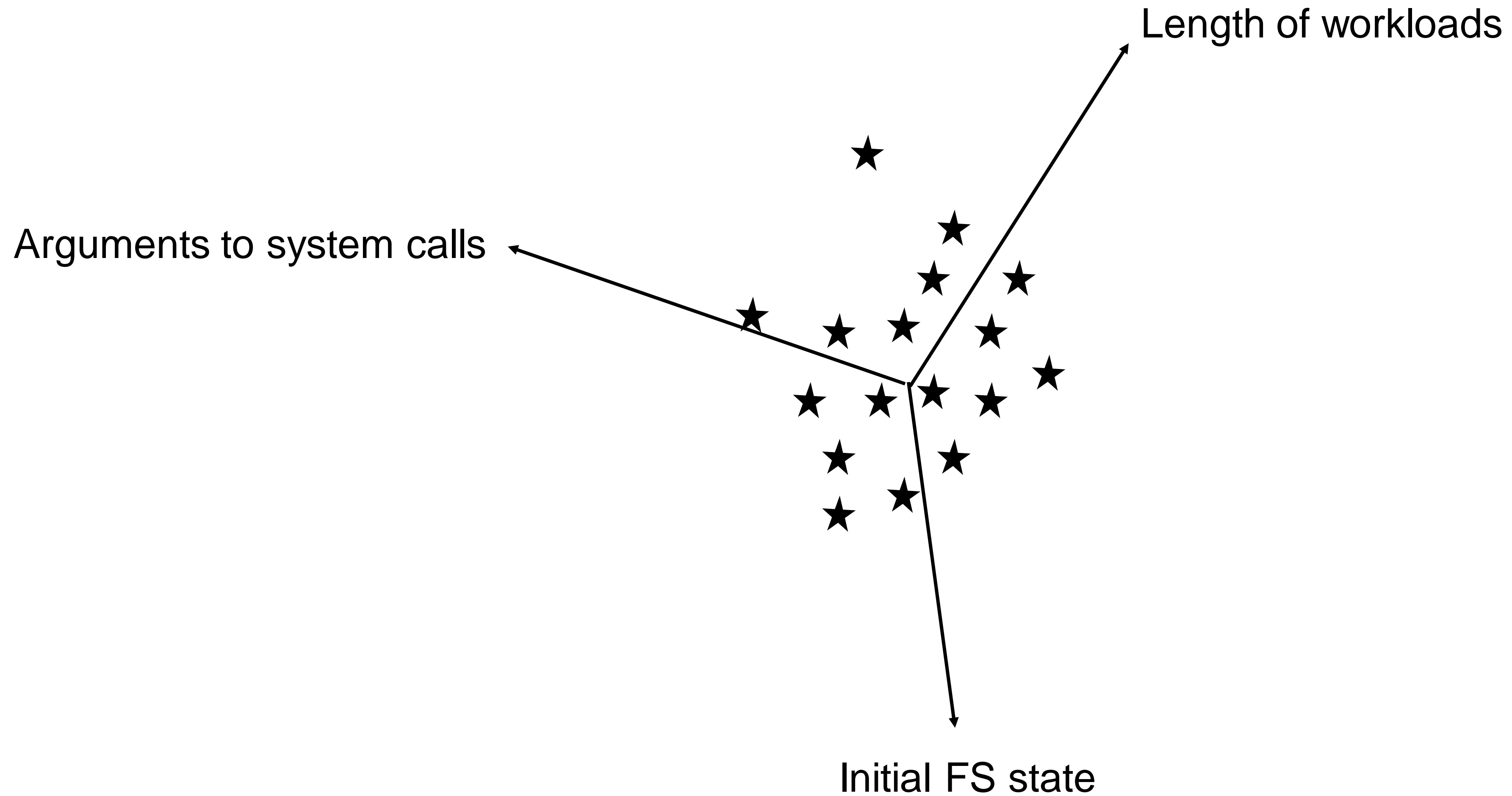
Outline

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

B³ : Bounded Black Box Crash Testing



B³ : Bounded Black Box Crash Testing



B³ : Bounded Black Box Crash Testing

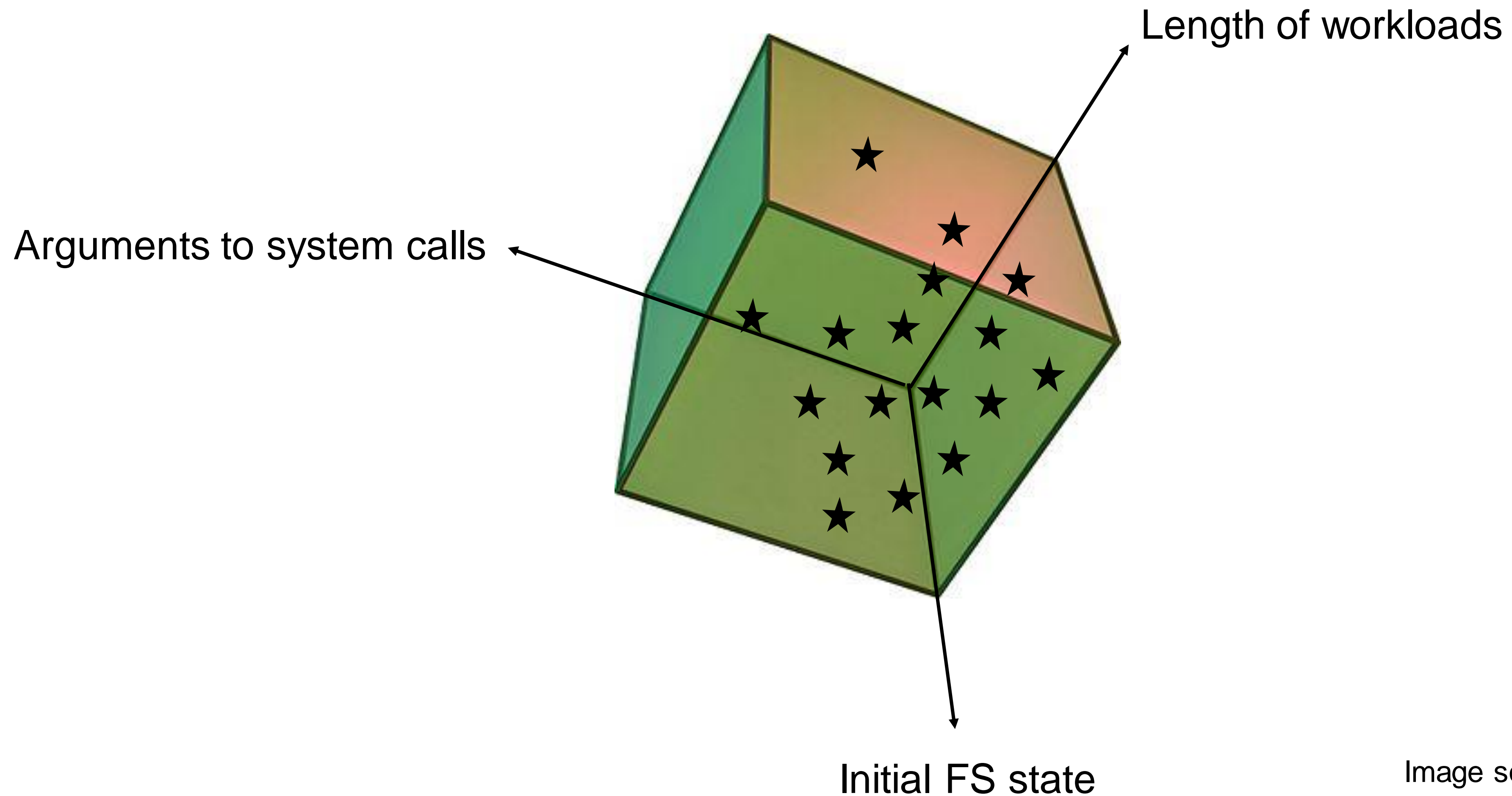
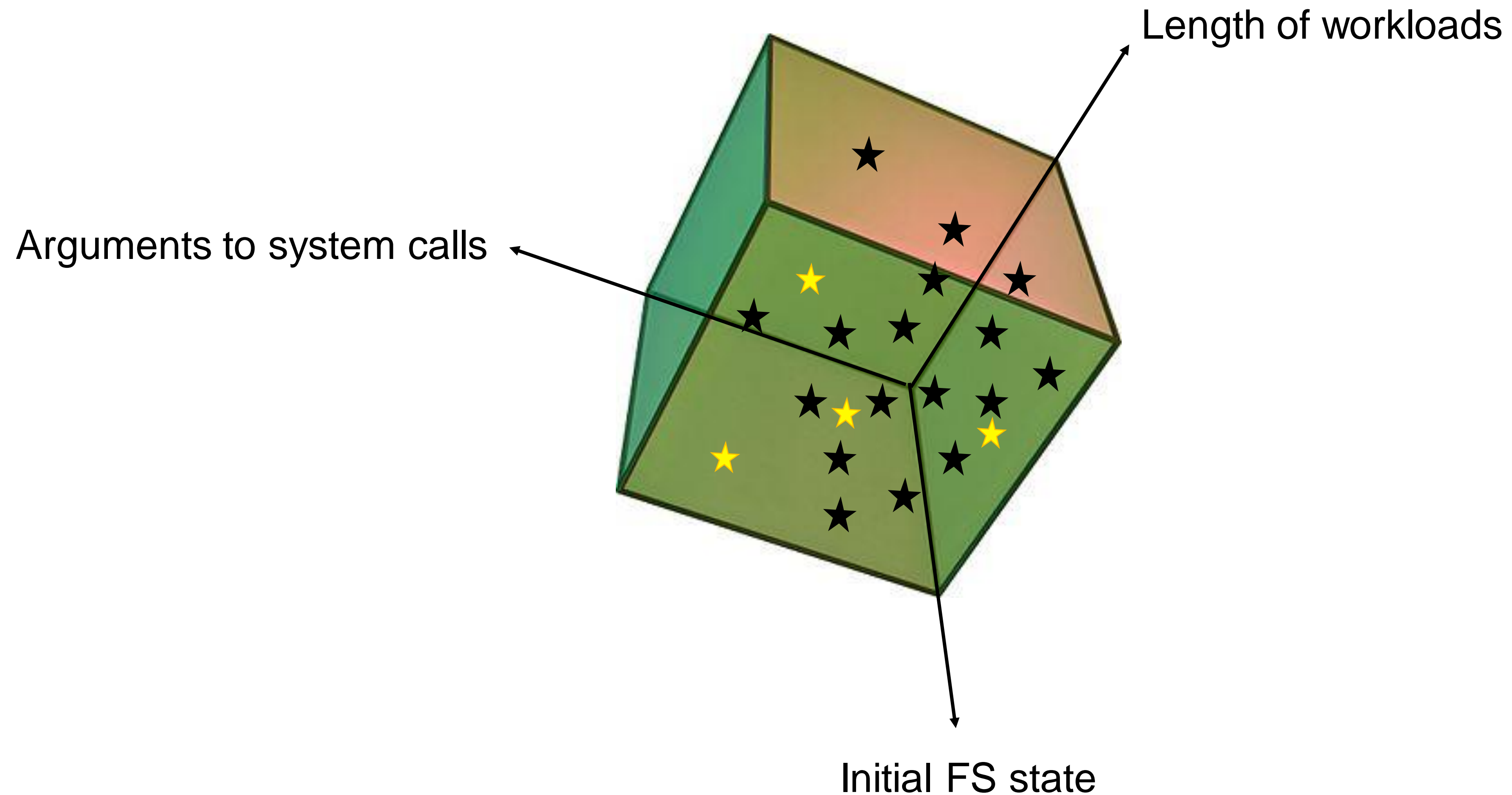
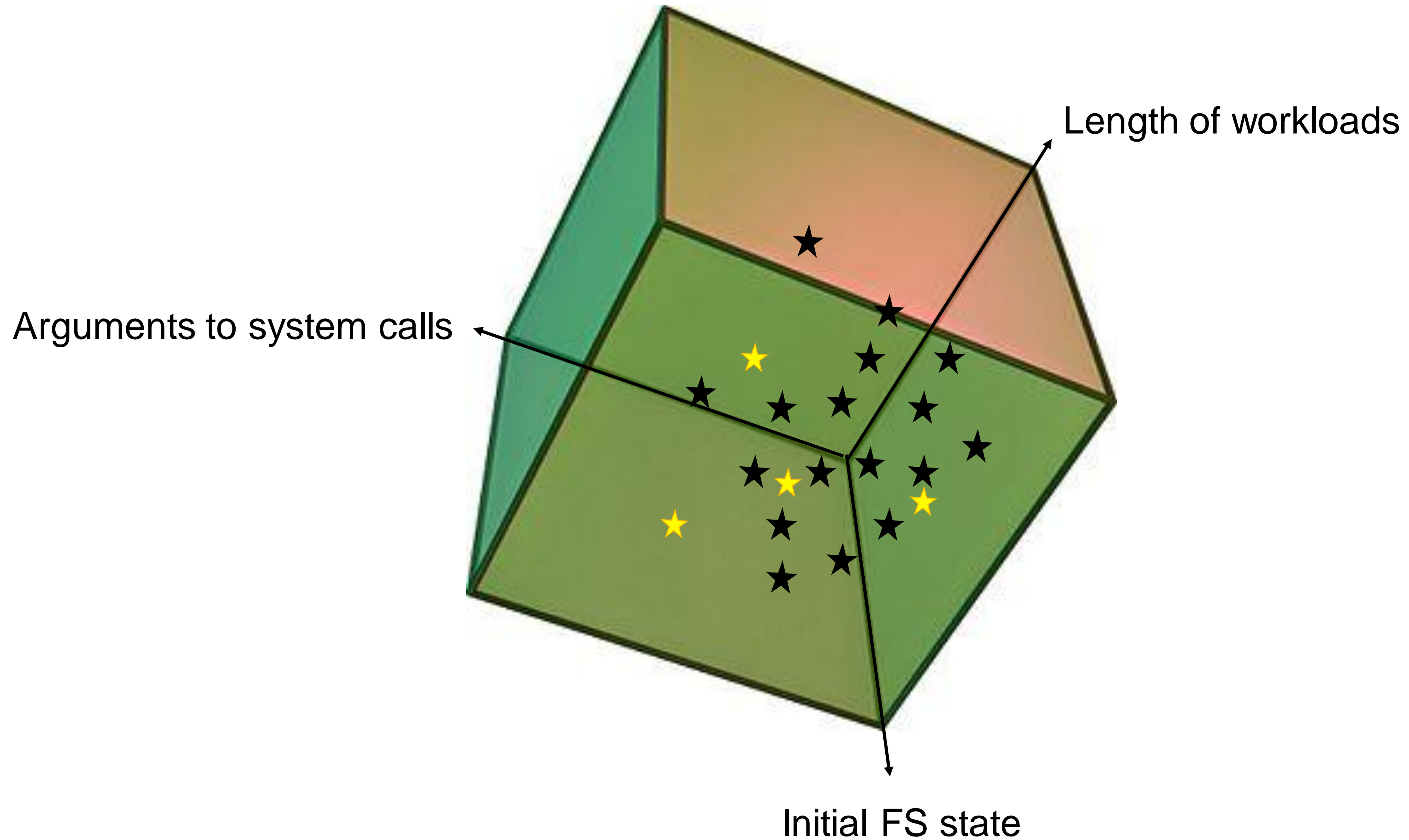


Image source: <https://en.wikipedia.org/wiki/Cube>

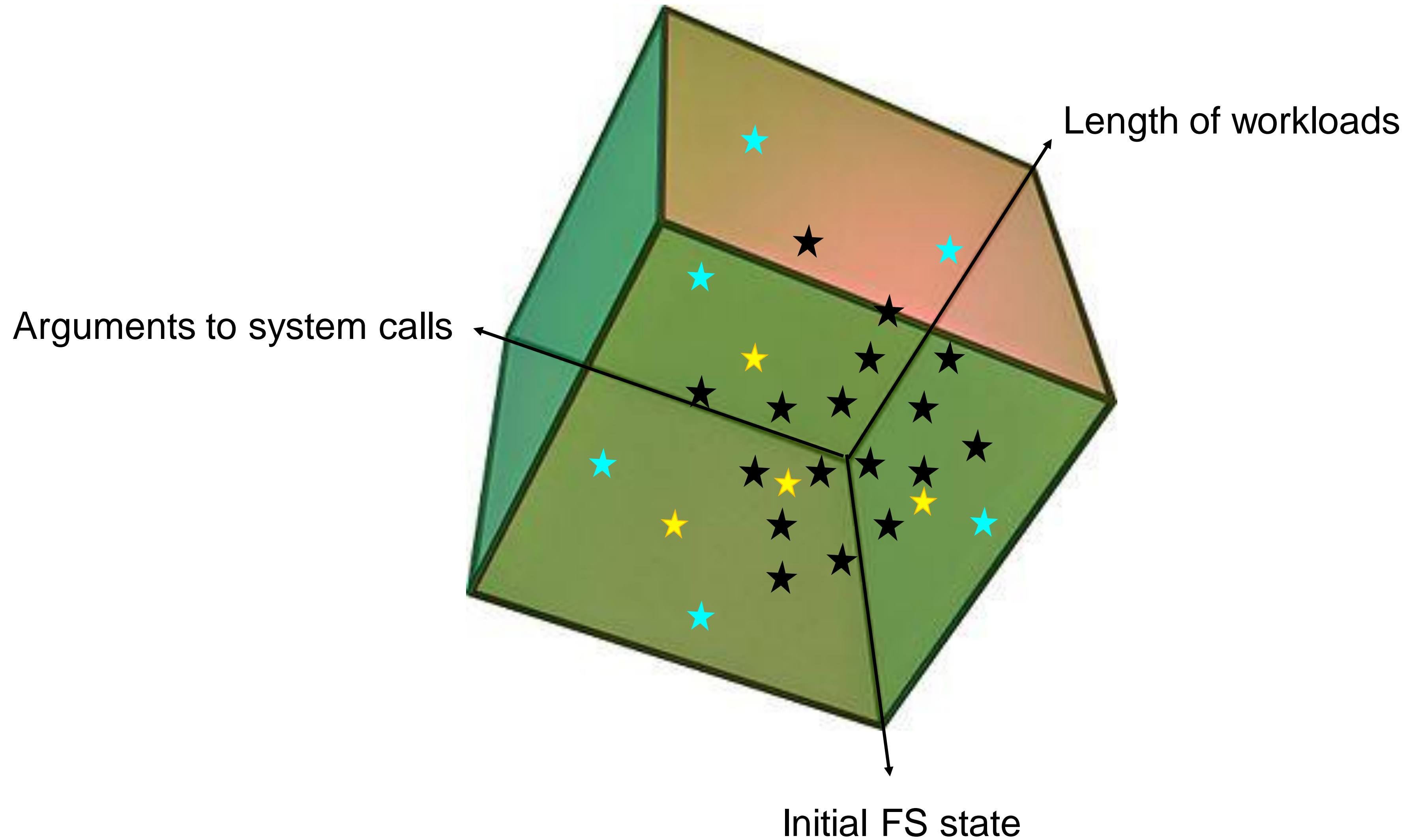
B³ : Bounded Black Box Crash Testing



B³ : Bounded Black Box Crash Testing



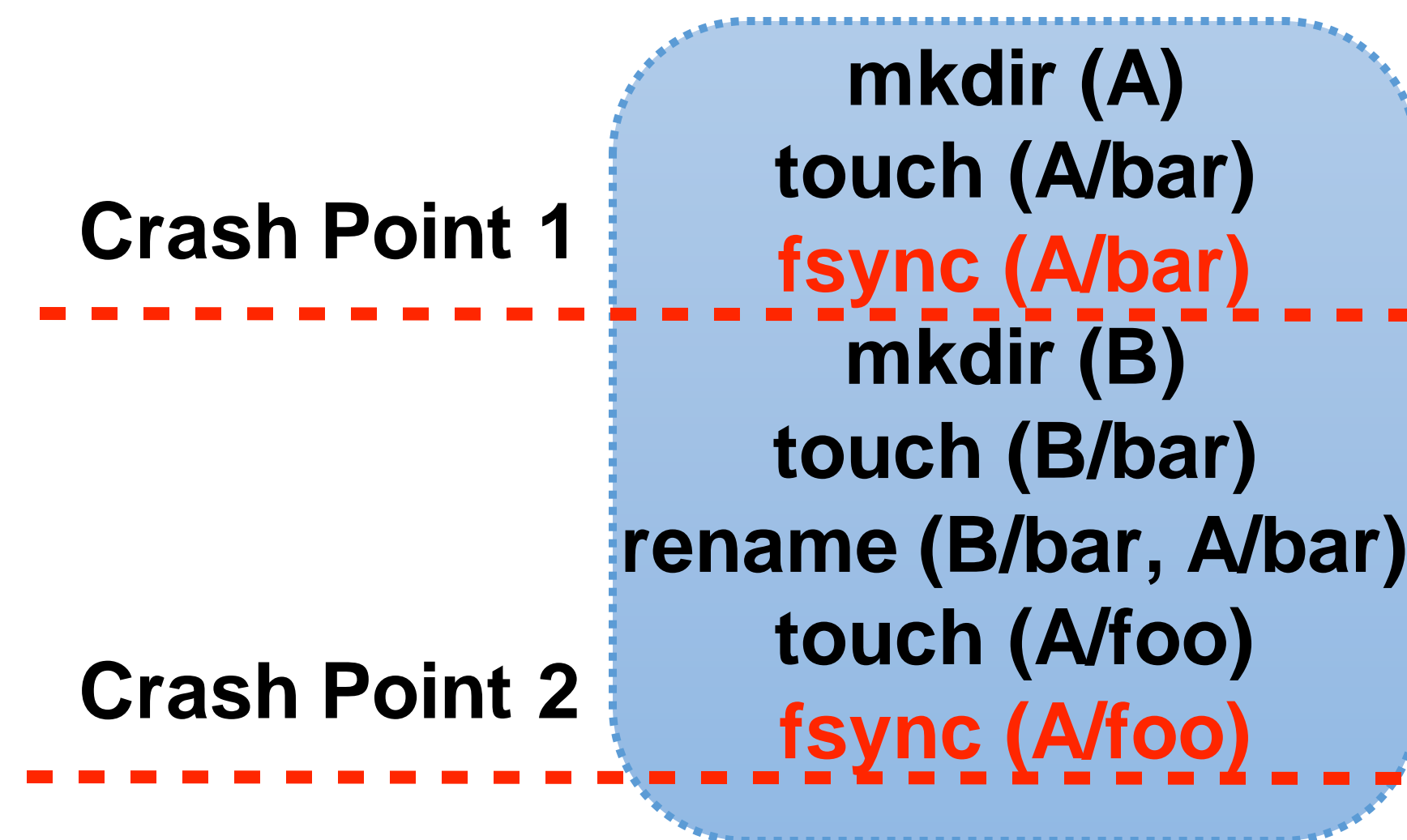
B³ : Bounded Black Box Crash Testing



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.

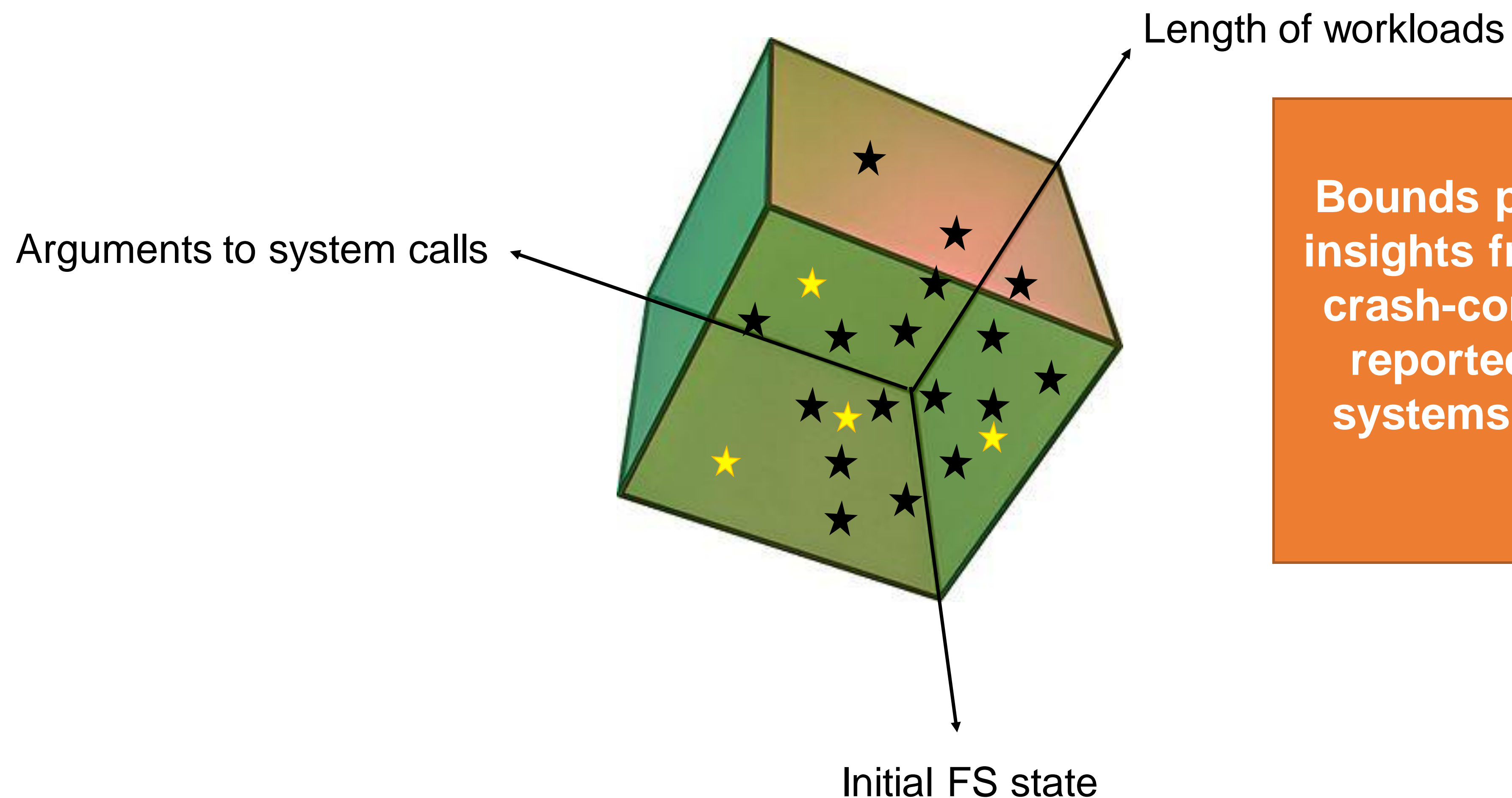
Limitations of B³

- 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

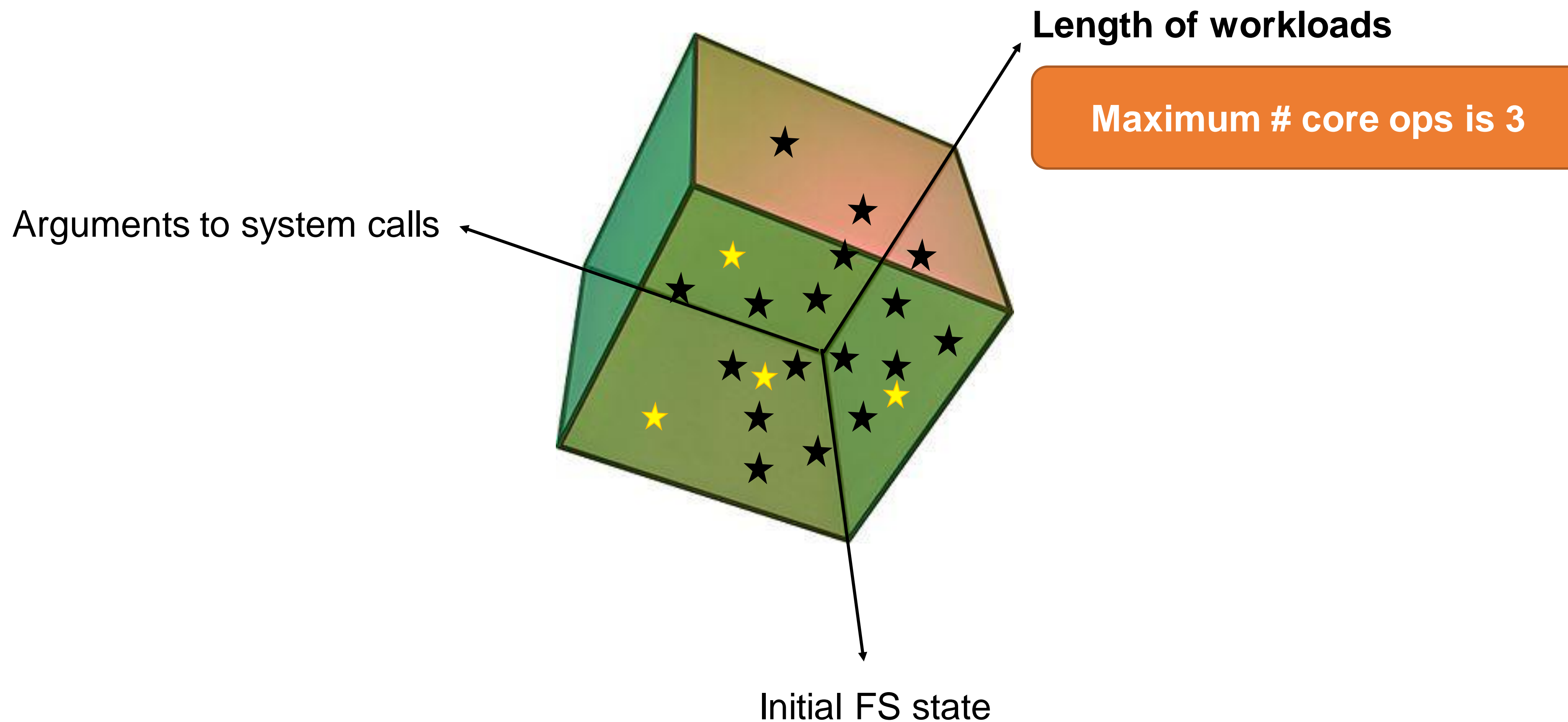
- CrashMonkey
- Bounded Black Box Crash Testing
- **Automatic Crash Explorer (ACE)**
- Results

Bounds chosen by ACE

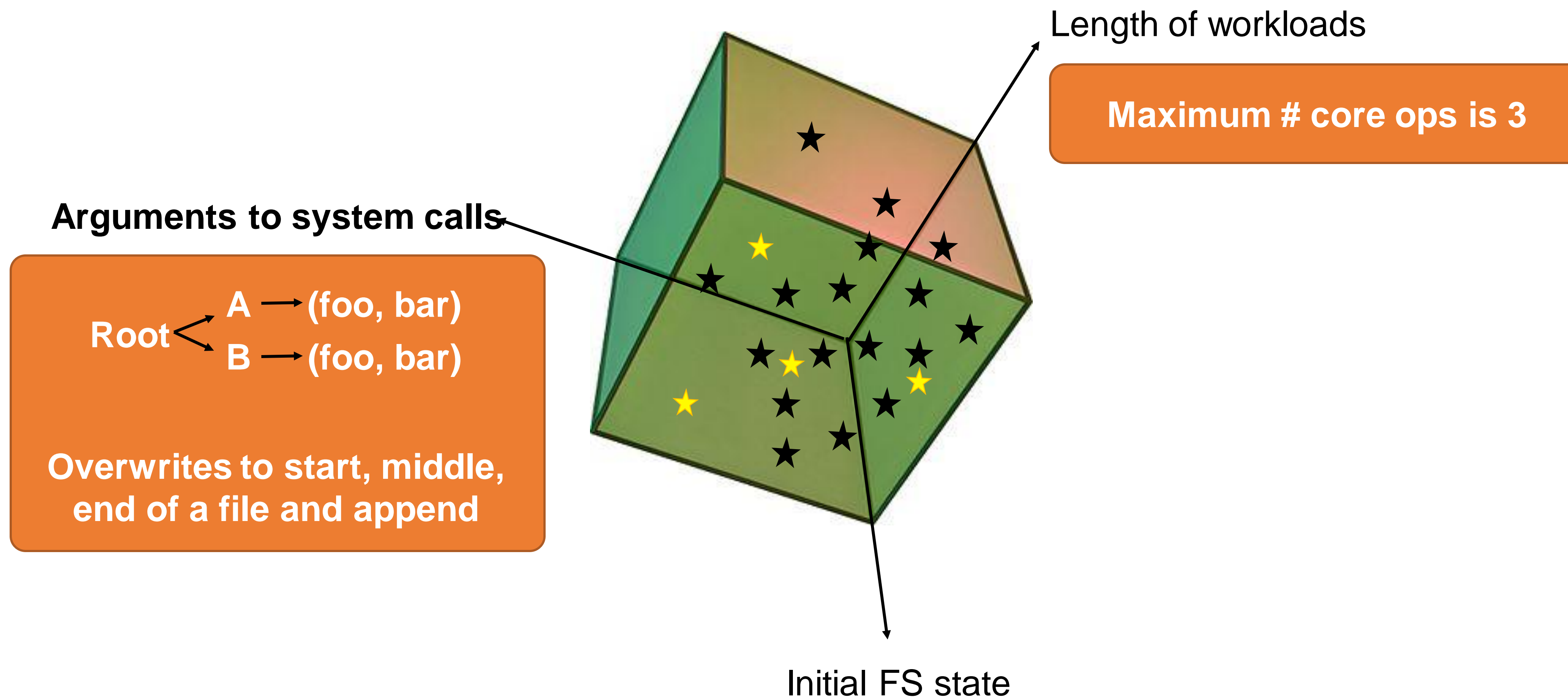


Bounds picked based on insights from the study of crash-consistency bugs reported on Linux file systems over the last 5 years

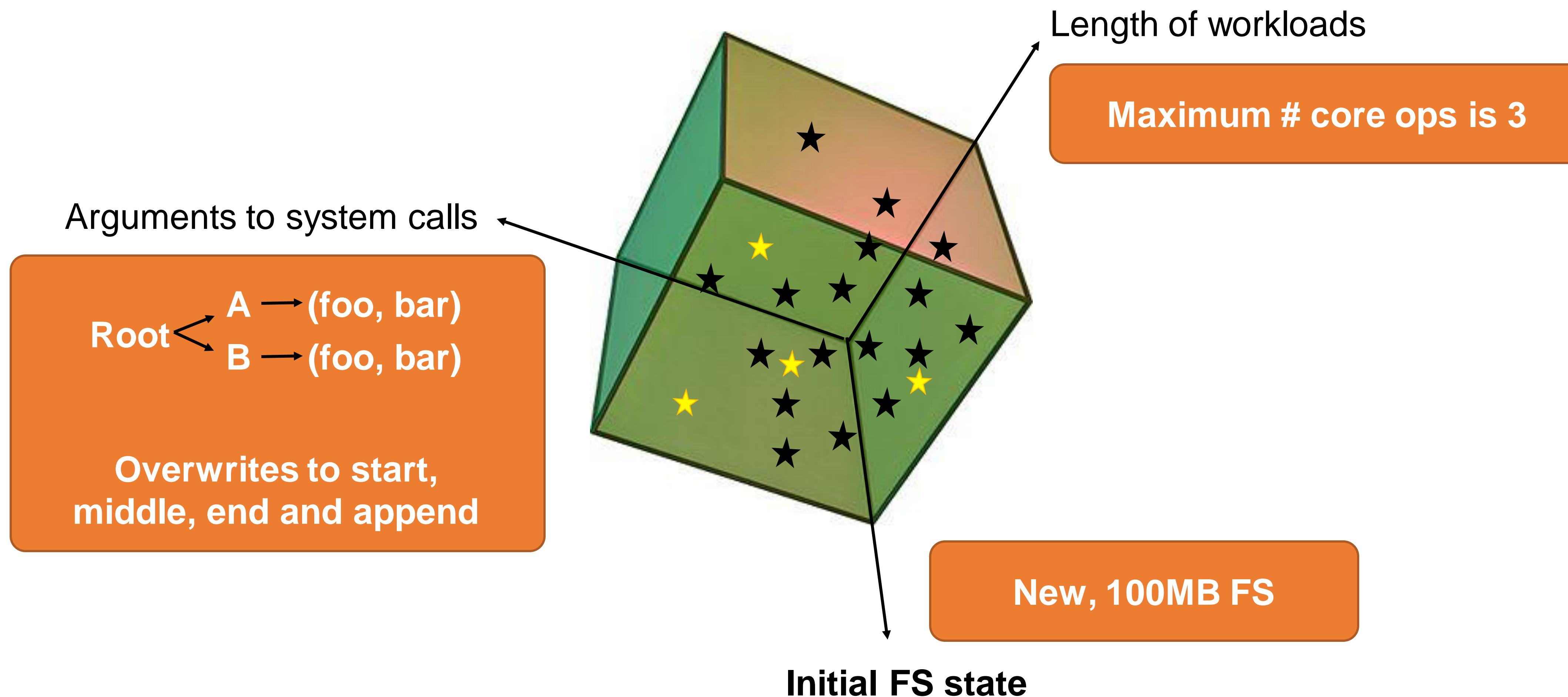
Bounds chosen by ACE



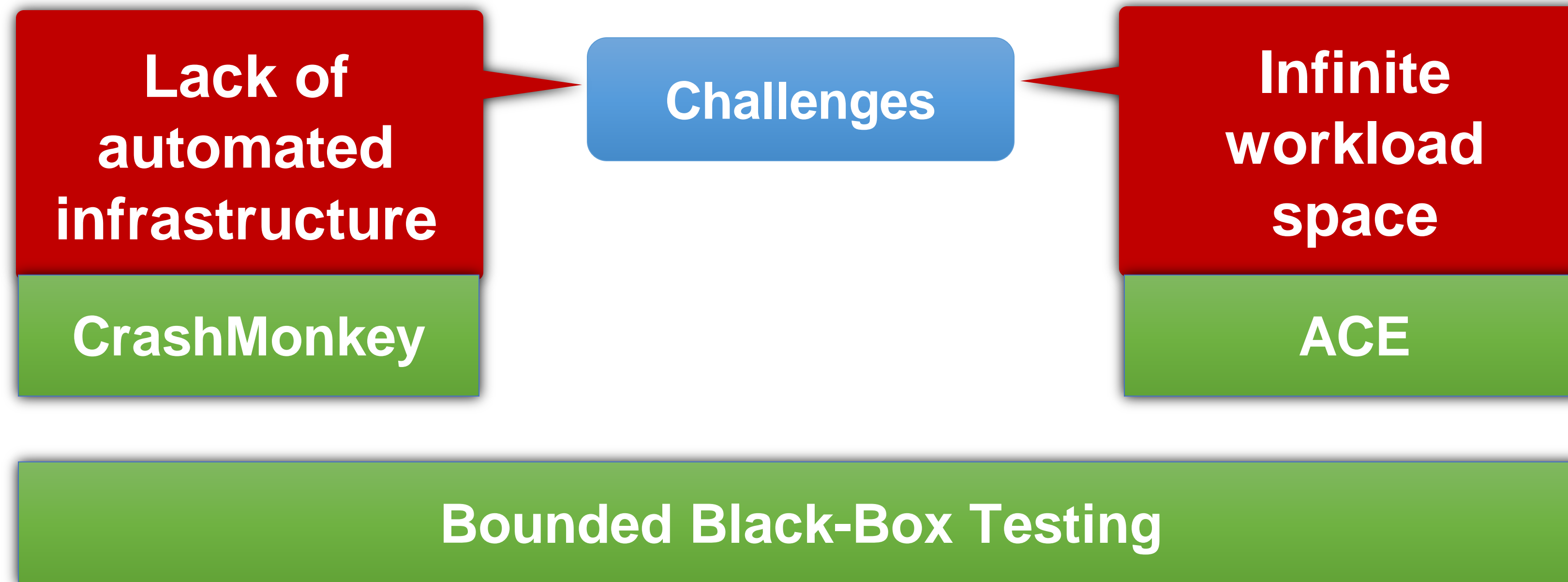
Bounds chosen by ACE



Bounds chosen by ACE



Challenges with Systematic Testing



Outline

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

Results

Sequence Length	# workloads	# Bugs Reproduced	# Bugs found
1, 2, 3	3.37M	24	10

- 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

Bounded Black-Box Crash Testing

- B^3 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 described as 'CrashMonkey: tools for testing file-system reliability (OSDI 18)'. It has 815 commits, 63 branches, 1 release, 1 environment, 9 contributors, and is licensed under Apache-2.0. The repository is currently on the 'master' branch. A recent merge pull request #129 from ashmrtn/port-4.14 is shown, with the latest commit d2c834e on May 28. The repository contains several folders: ace, code, docs, googletest @ aa148eb, setup, test, and vm_scripts.

Folder	Description	Last Commit
ace	Fixing reported issues in Ace	8 months ago
code	Fix comment about struct field kernel version.	7 months ago
docs	Merge pull request #129 from ashmrtn/port-4.14	4 months ago
googletest @ aa148eb	Pull in gtests for C++ code.	2 years ago
setup	Add script for building VMs. Update README a tad more.	2 years ago
test	Add tests for unaligned <4k writes for WriteData.	last year
vm_scripts	Almost done with documentation	last year

Thanks!
Questions?

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