# Physical Disentanglement in a Container-Based File System

Lanyue Lu, Yupu Zhang, Thanh Do, Samer Al-Kiswany
Andrea C. Arpaci-Dusseau, Remzi H. Arpaci-Dusseau
University of Wisconsin Madison

# Motivation

Isolation in various subsystems

- ➤ resources: virtual machines, Linux containers
- ➤ security: BSD jail, sandbox
- ➤ reliability: address spaces

File systems lack isolation

- ➤ physical entanglement in modern file systems

Three problems due to entanglement

- ➤ global failures, slow recovery, bundled performance

# Global Failures

## Definition

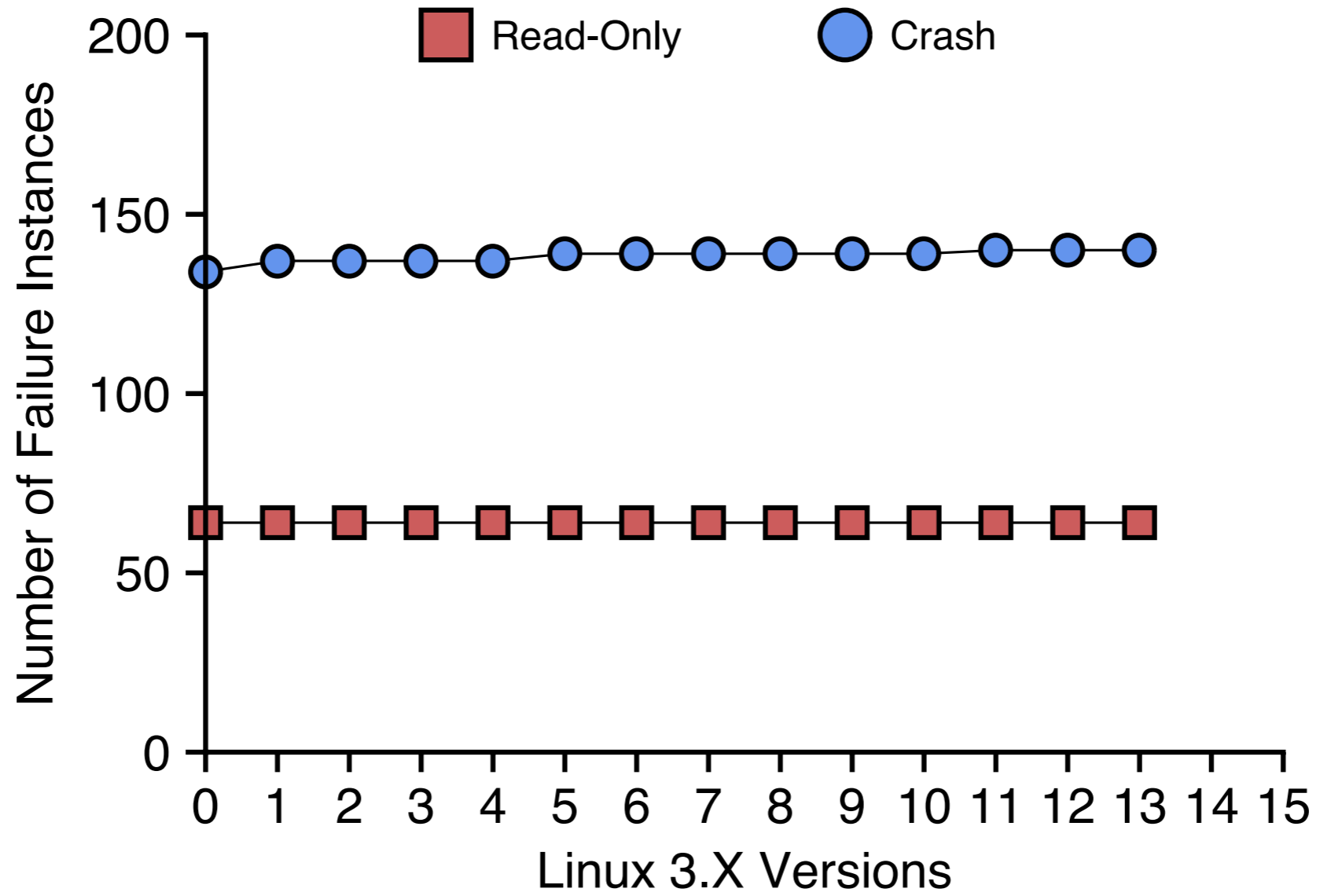- a failure which impacts all users of the file system or even the operating system

## Read-Only

- mark the file system as read-only
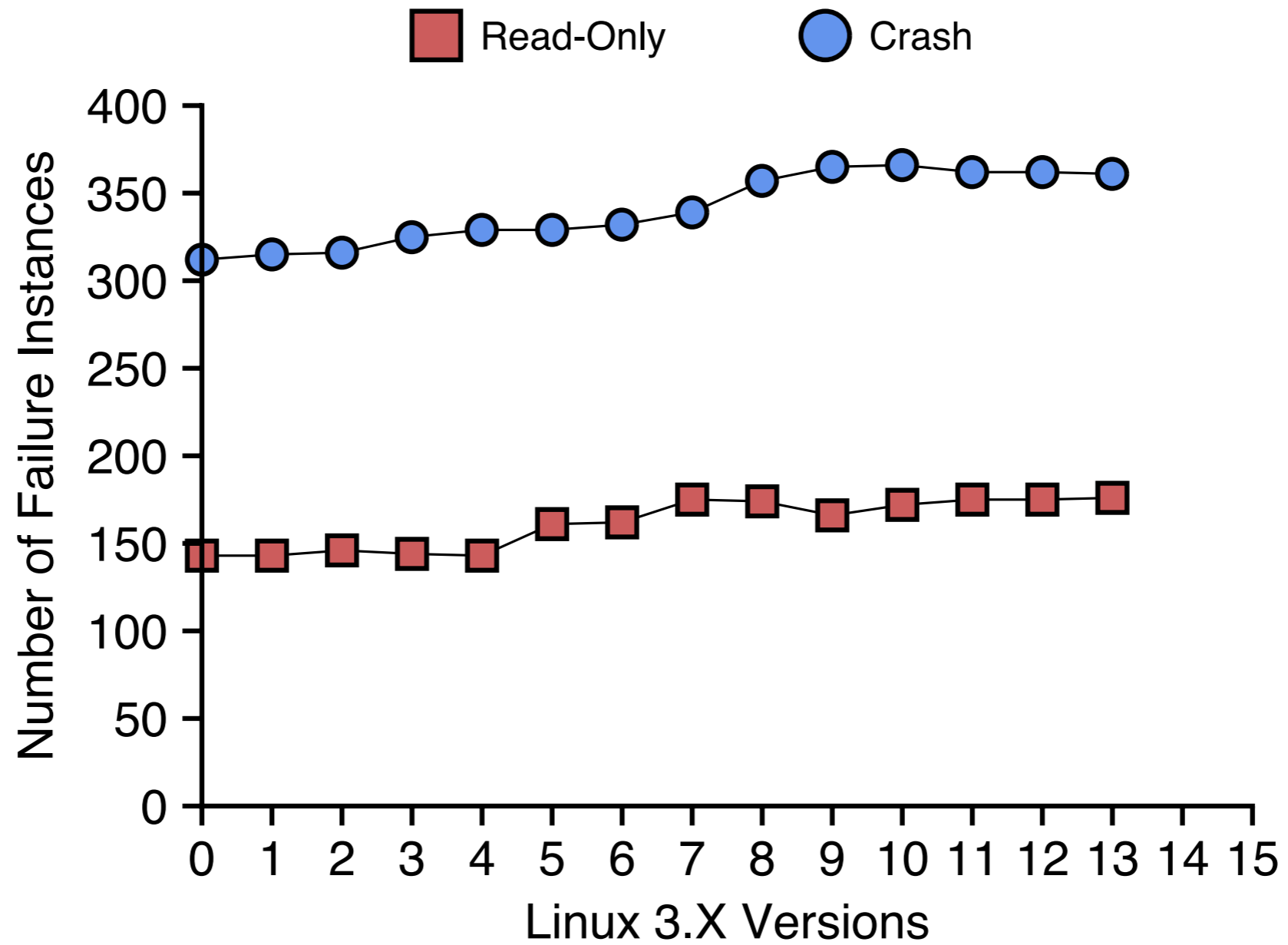- e.g., metadata corruption, I/O failure

## Crash

- crash the file system or the operating system
- e.g., unexpected states, pointer faults
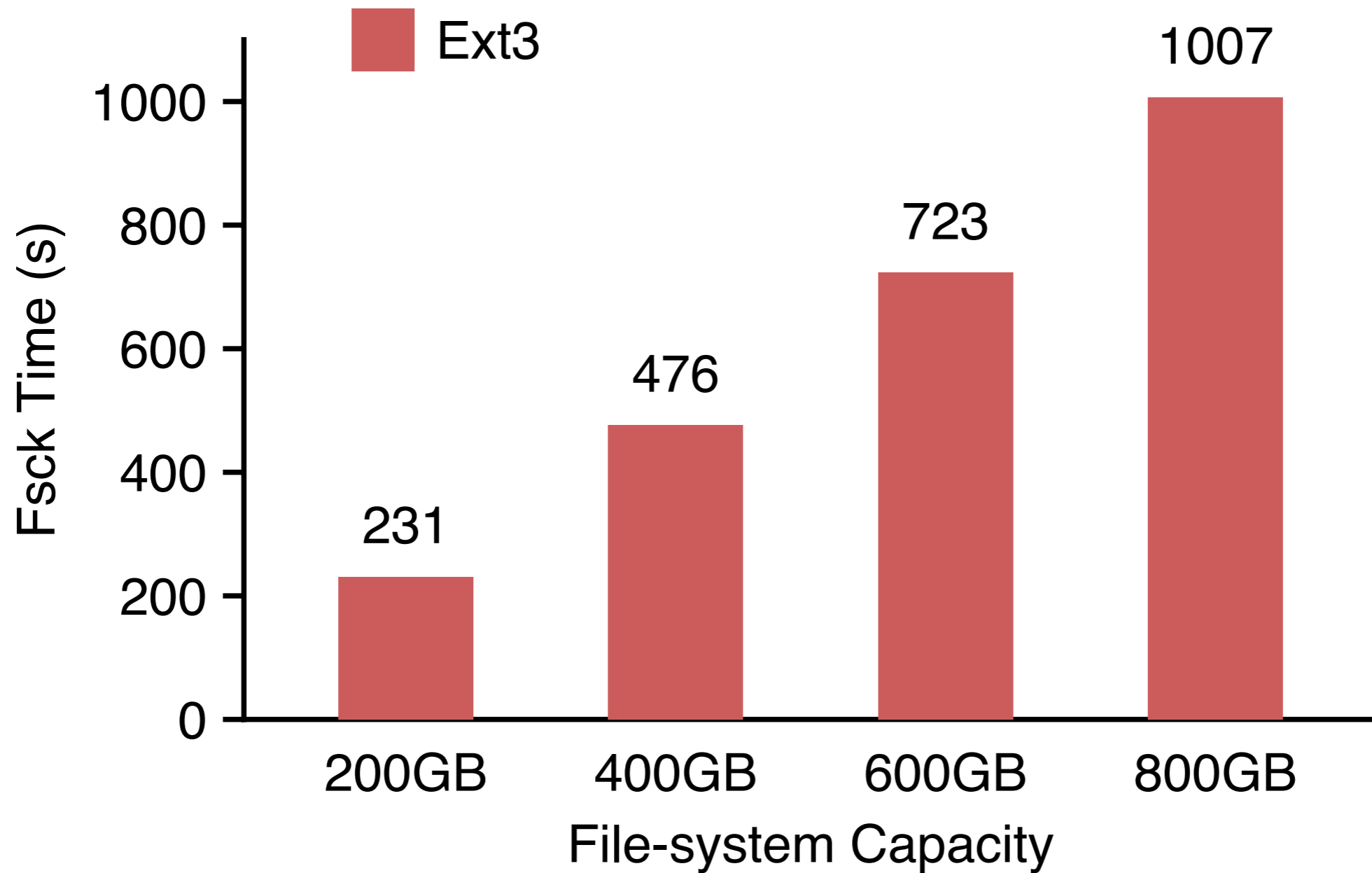
# Ext3

# Ext4

# Slow Recovery

## File system checker
- repair a corrupted file system
- usually offline

## Current checkers are not scalable
- increasing disk capacities and file system sizes
- scan the whole file system
- checking time is unacceptably long

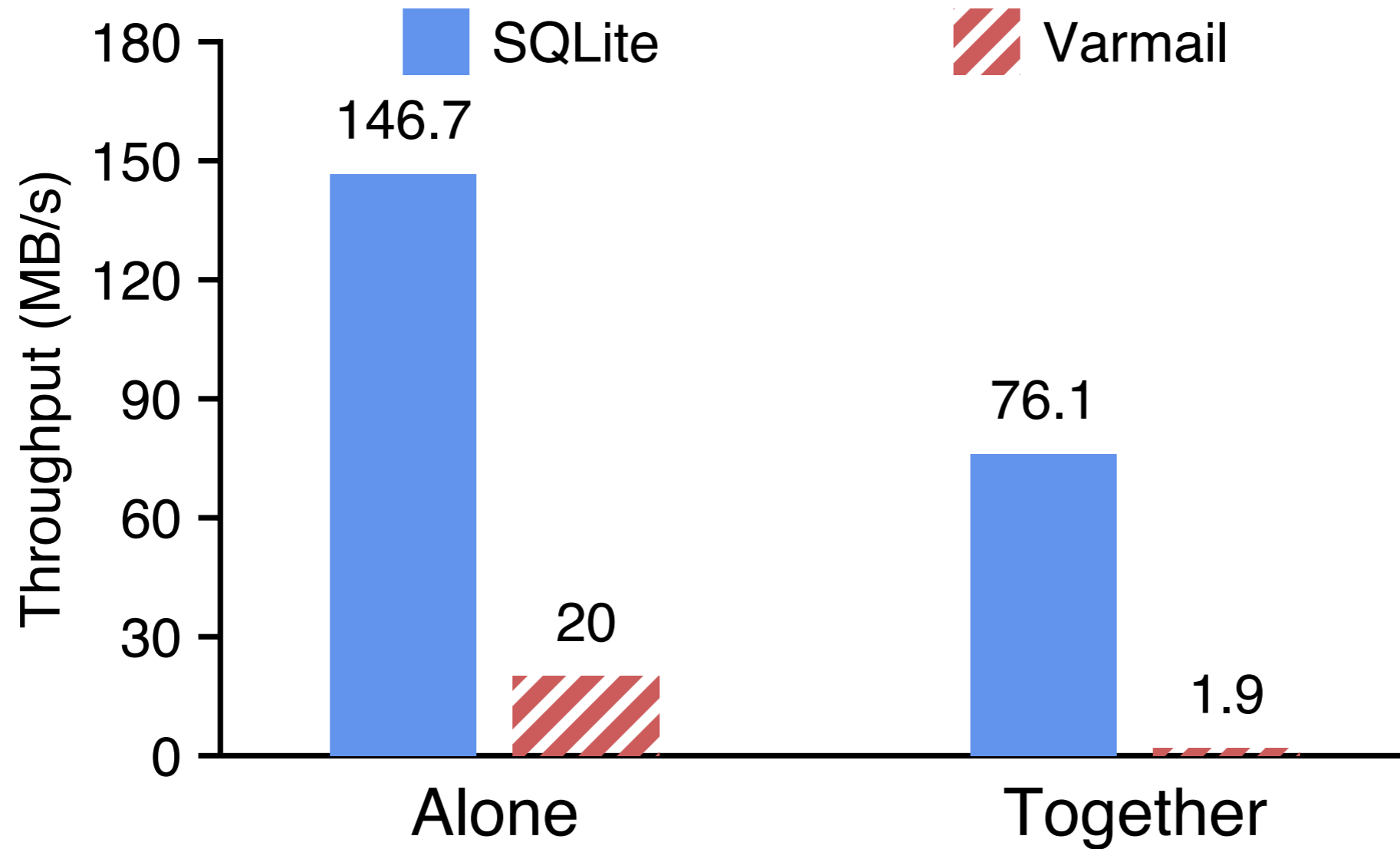# Scalability of fsck on Ext3

# Bundled Performance

Shared transaction
- all updates share a single transaction
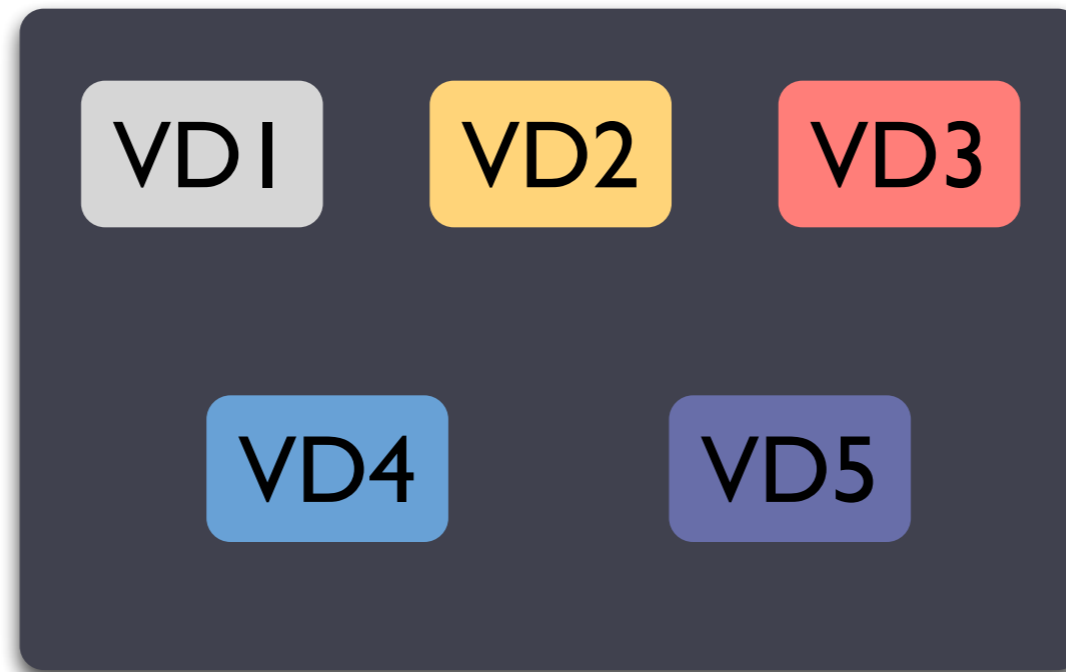- unrelated workloads affect each other

Consistency guarantee
- the same journal mode for all files
- limited flexibility for different tradeoffs

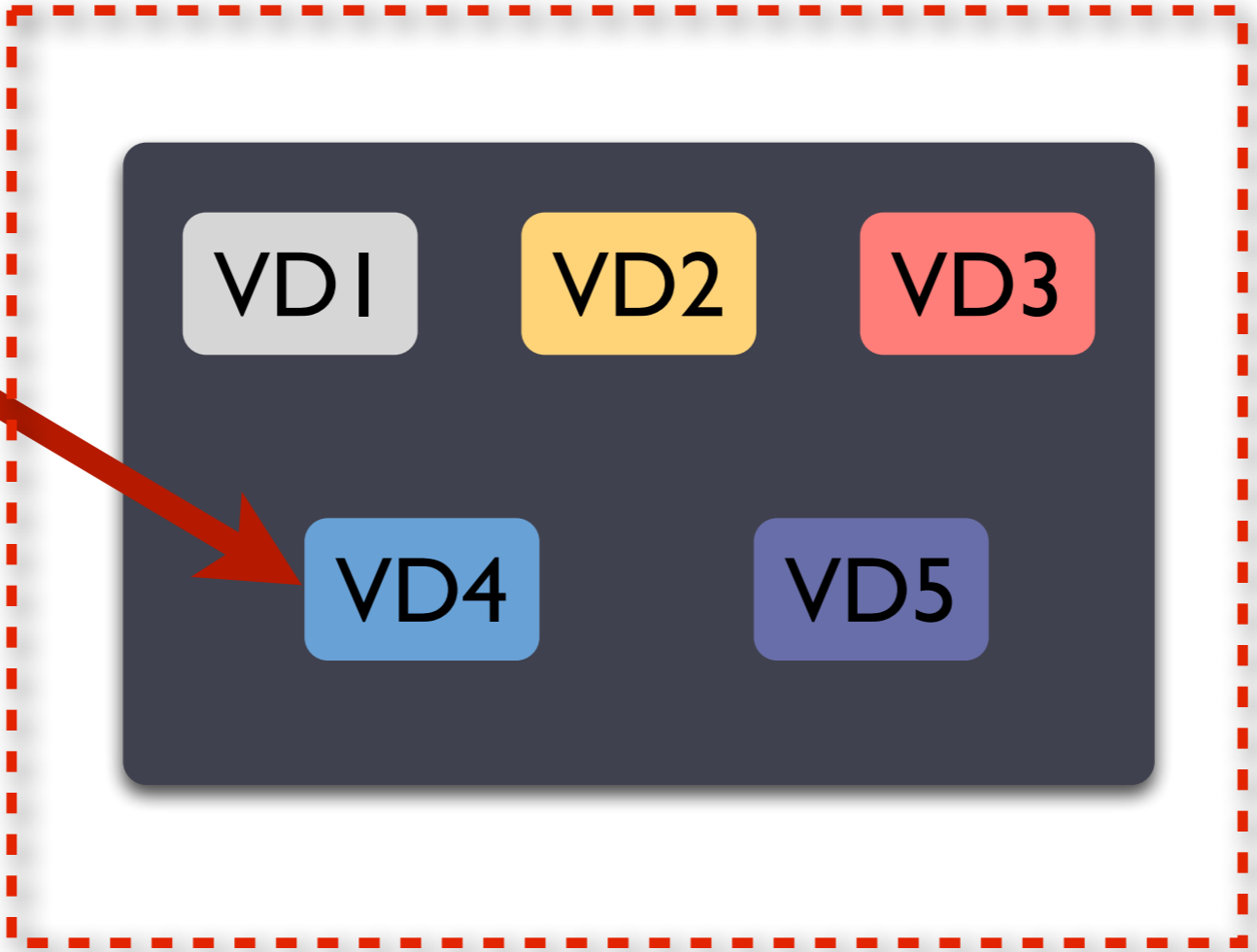# Bundle Performance on Ext3

# Server Virtualization

VM1 VM2 VM3 VM4 VM5

VD1 VD2 VD3
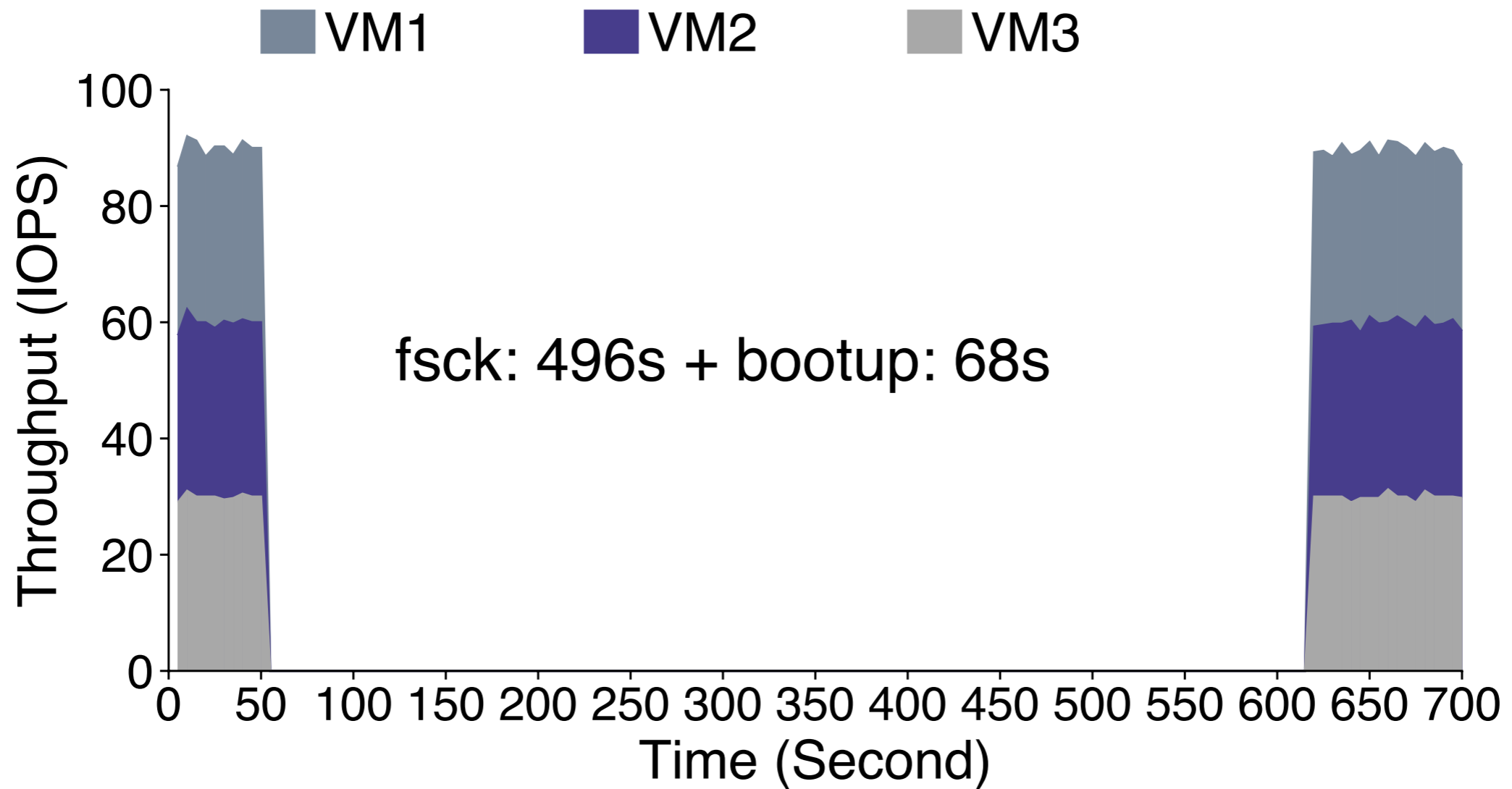
VD4 VD5

Shared file system

# Virtual Machines



fsck: 496s + bootup: 68s

# Our Solution

A new abstraction for disentanglement

A disentangled file system: IceFS

Three major benefits of IceFS
- failure isolation
- localized recovery
- specialized journaling
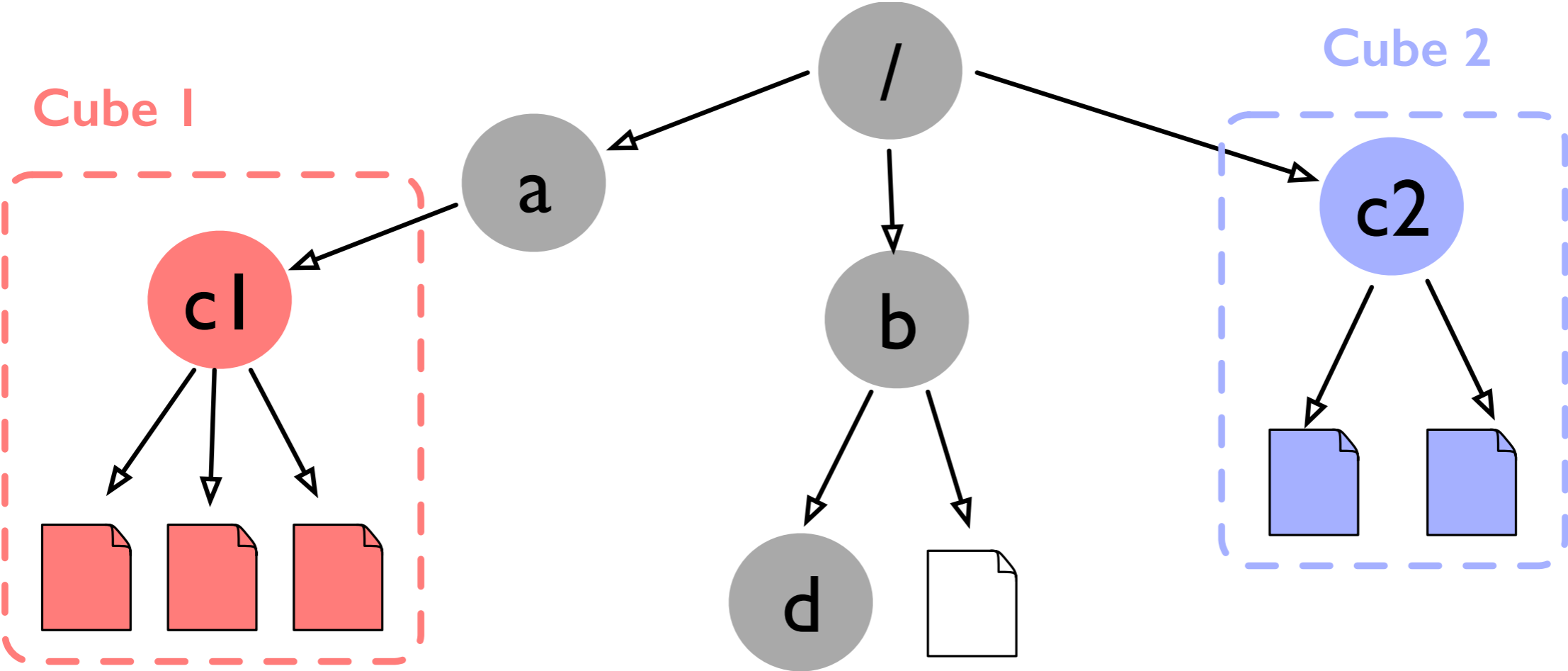
# A New Abstraction: Cube

## What is a cube
- an independent container for a group of files
- physically isolated

## Interface
- create a cube: mkdir(cube_flag)
- delete a cube: rmdir()
- add files to a cube
- remove files to a cube
- set cube attributes

# Cube Example

# A Disentangled File System

## No shared physical resources
- ➥ no shared metadata: e.g., block groups
- ➥ no shared disk blocks or buffers

## No access dependency
- ➥ partition linked lists or trees
- ➥ avoid directory hierarchy dependency

## No bundled transactions
- ➥ use separate transactions
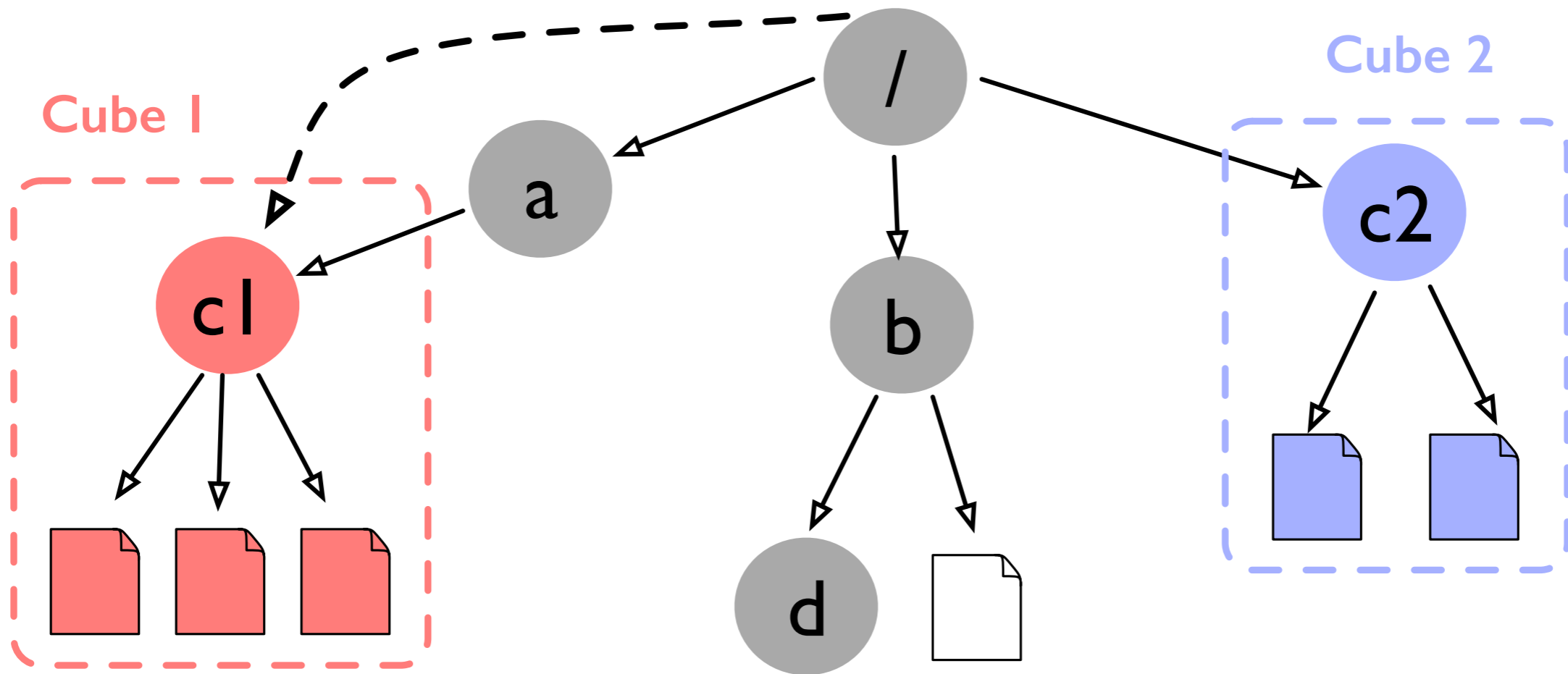- ➥ enable customized journaling modes

# IceFS

A prototype based on Ext3 in Linux 3.5

Disentanglement
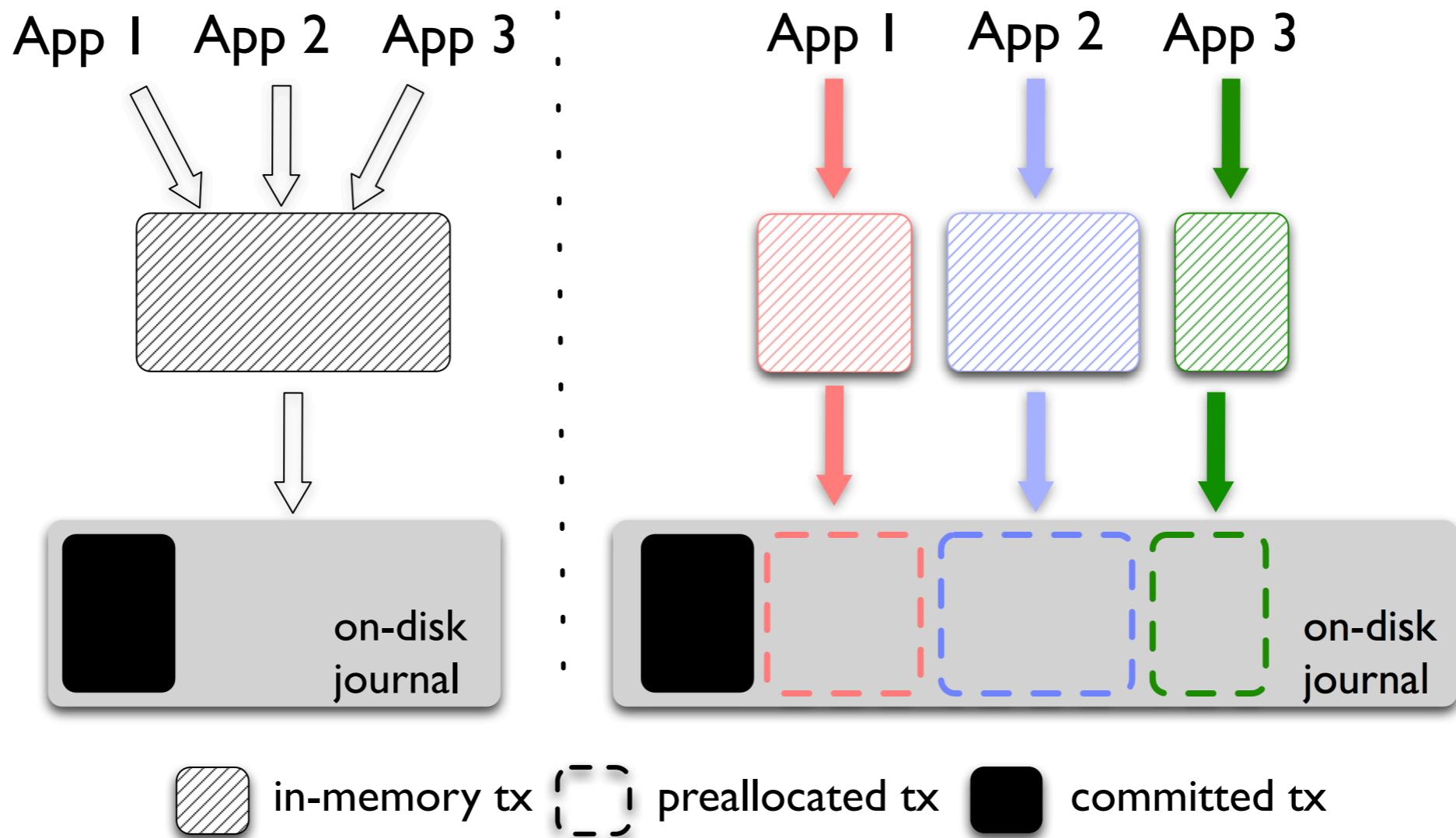- directory indirection
- transaction splitting

# Directory Indirection



Pathname matching for cubes' top directory paths

Cubes' dentries are pinned in memory

# Transaction Splitting



(a) Ext3/Ext4  (b) IceFS

Legend: in-memory tx, preallocated tx, committed tx

Separate transactions from different cubes commit to the disk journal in parallel

# Benefits of Disentanglement

## Localized reactions to failures
- ➡ per-cube read-only and crash

## Localized recovery
- ➡ only check faulty cubes
- ➡ offline and online

## Specialized journaling
- ➡ parallel journaling
- ➡ diverse journal modes
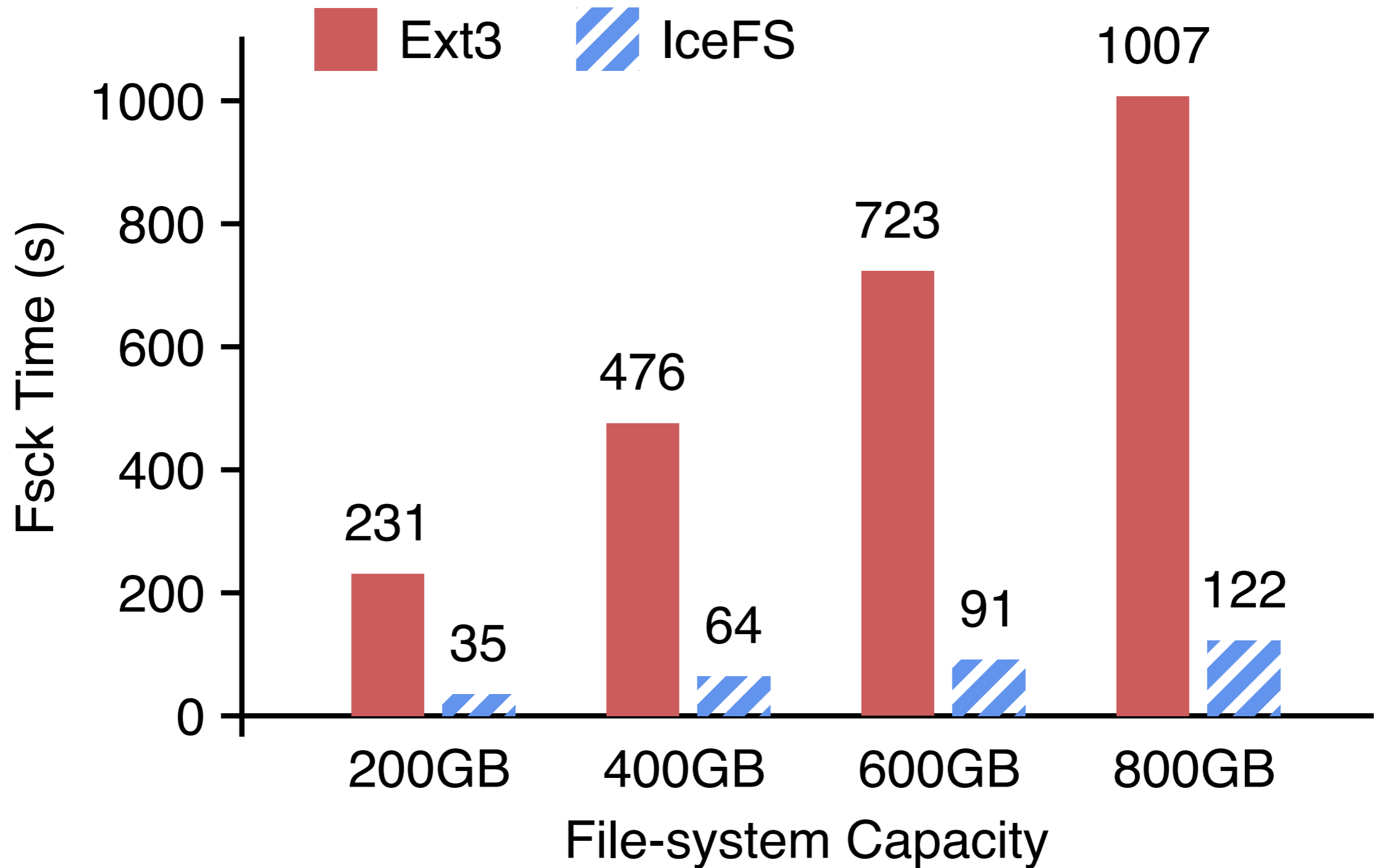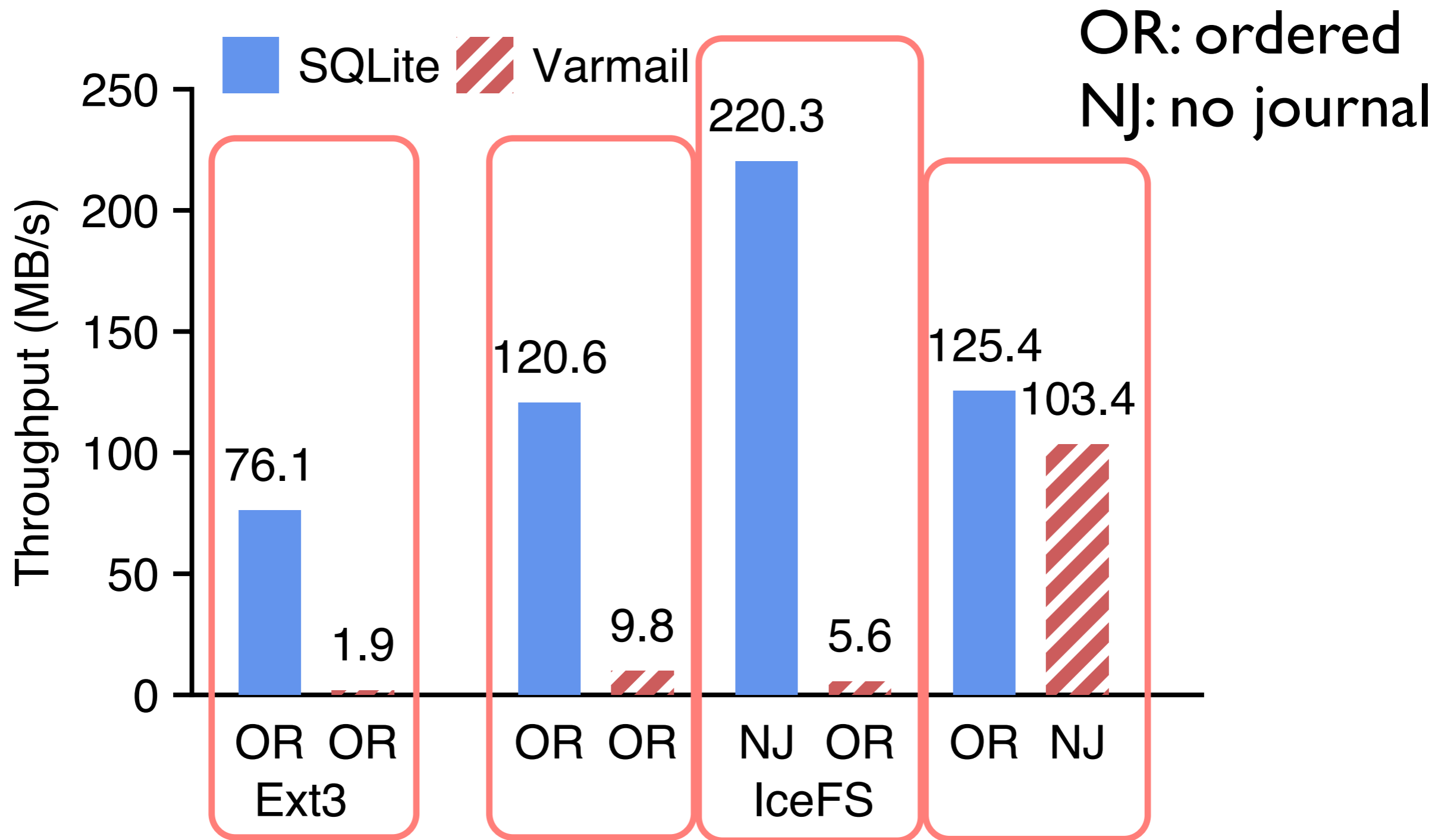  - ➡ no journal, no fsync, writeback, ordered, data

Motivation

IceFS
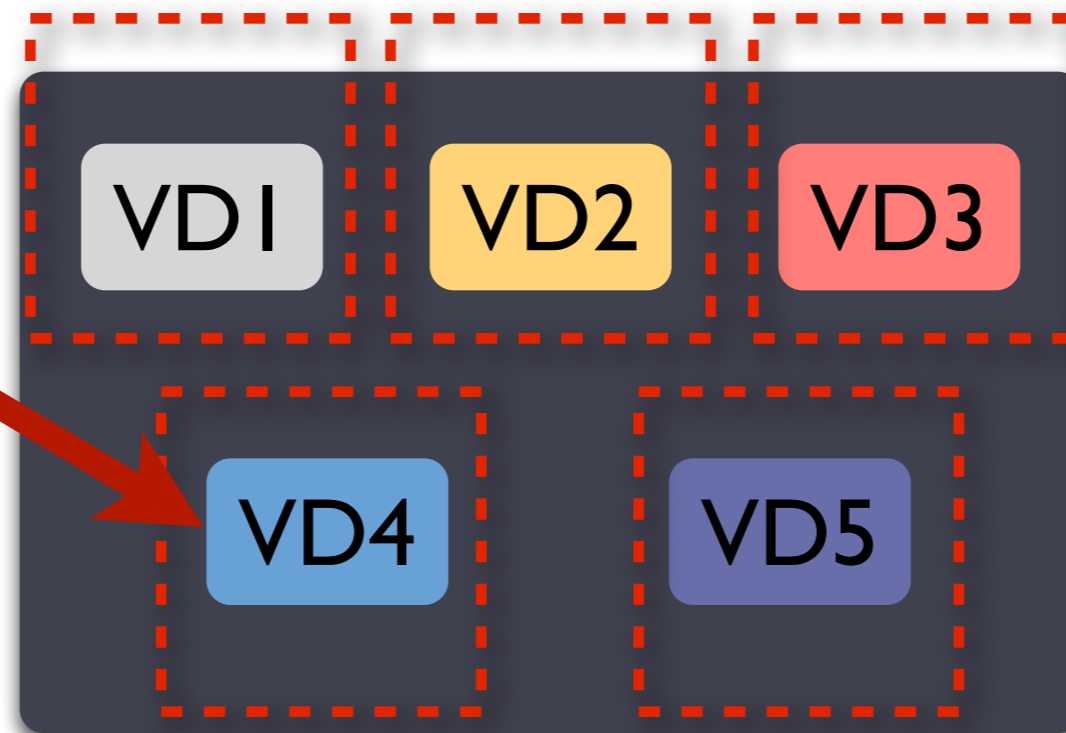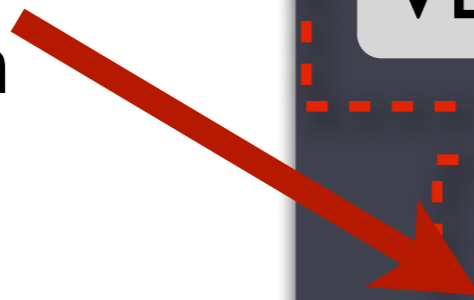
Evaluation

Conclusion

# Fast Recovery

# Server Virtualization



VM1　VM2　VM3　VM4　VM5

metadata corruption

VD1　VD2　VD3

VD4　VD5
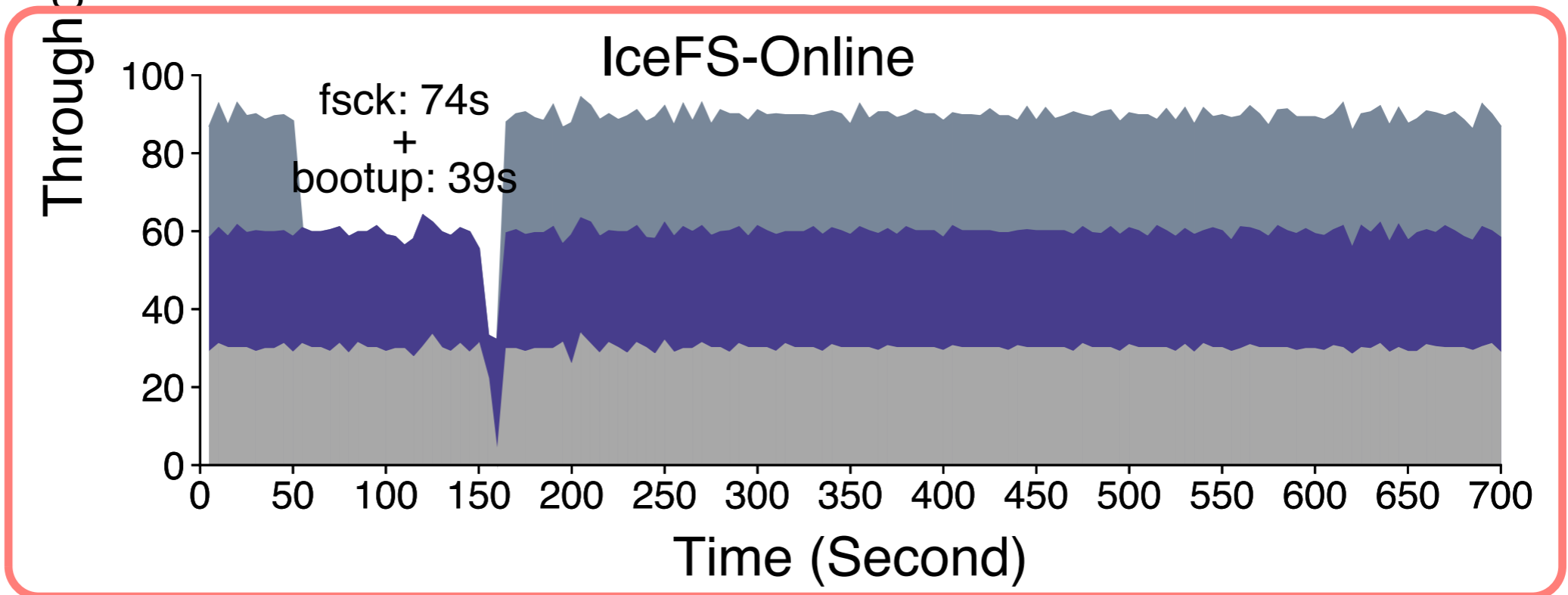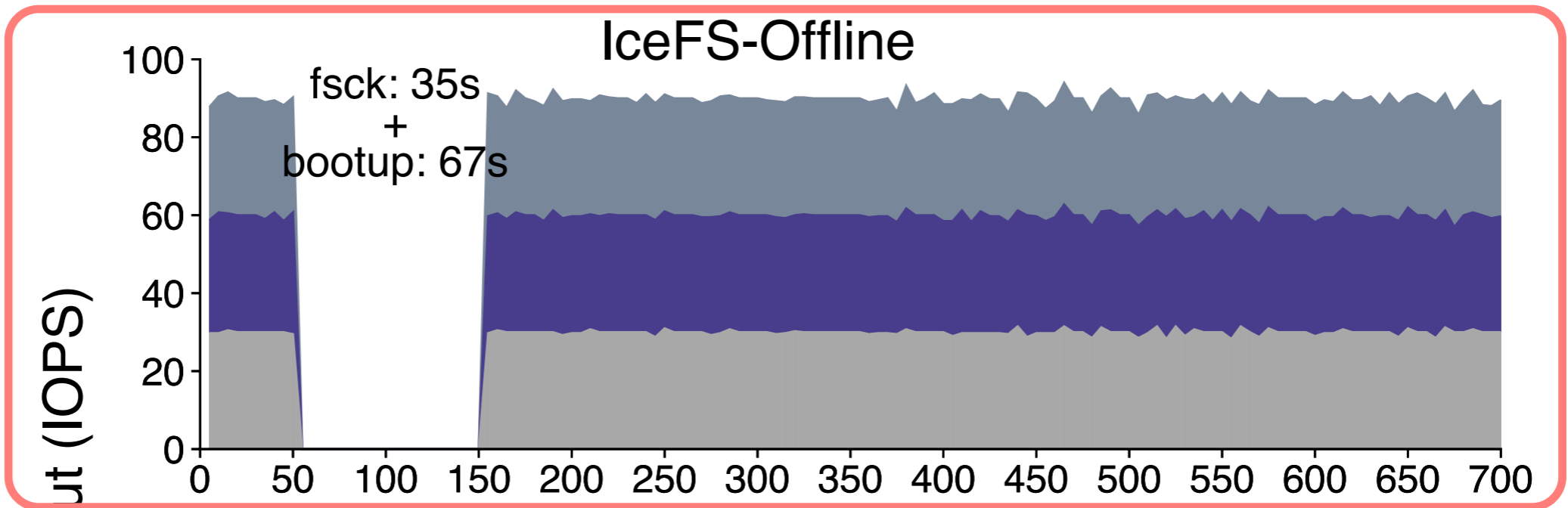
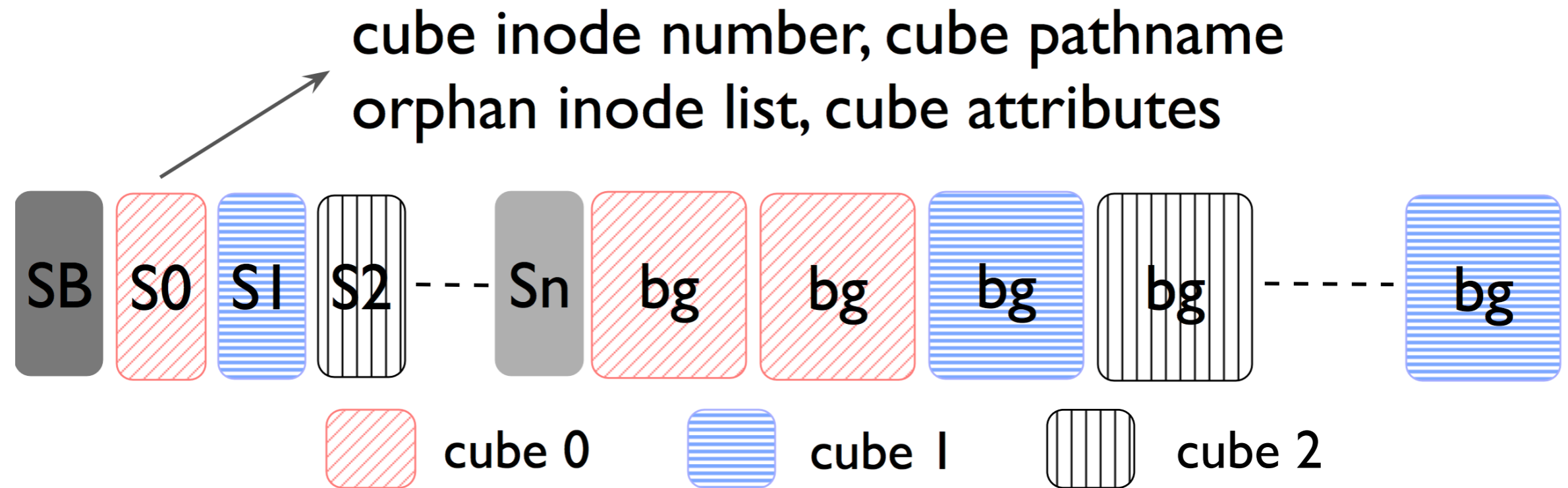Shared file system

# Server Virtualization

# Conclusion

File systems lack physical isolation

Our contribution
- a new abstraction: cube
- a disentangled file system: IceFS
- demonstrate its benefits:
  - isolated failures
  - localized recovery
  - specialized journaling
  - improving both reliability and performance

# Questions ?

# IceFS Disk Layout



cube inode number, cube pathname
orphan inode list, cube attributes

SB  S0  S1  S2  ----  Sn  bg  bg  bg  bg  ----  bg

cube 0    cube 1    cube 2

Each cube has one sub-super block (Si) and its own
block groups (bg)