# First Responder: Persistent Memory Simultaneously as High Performance Buffer Cache and Storage

**Hyunsub Song**

**Shean Kim**  **J. Hyun Kim**  **Ethan JH Park**  **Sam H. Noh**

**USENIX ATC 2021**

**UNIST**

**(Ulsan National Institute of Science and Technology)**

UNIST
ULSAN NATIONAL INSTITUTE OF SCIENCE AND TECHNOLOGY
2009

ECE
SCHOOL OF ELECTRICAL AND COMPUTER ENGINEERING

CISSR
Center for Innovative System Software Research

NECSST
Next-generation Embedded / Computer System Software Technology

# TABLE OF CONTENTS

NECSST
Next-generation Embedded / Computer System Software Technology

# TABLE OF CONTENTS

Next-generation Embedded / Computer System Software Technology

**Persistent Memory Features**

- Non-volatility
- Byte-level random access
- Fast access time (*nanoseconds*)

**Storage** + **Memory**



Persistent memory is evolving

HDD (~10ms)    SAS SSD (~150µs)    NVMe SSD (<100µs)

**NECSST**
Next-generation Embedded / Computer System Software Technology

| | | | |
|---|---|---|---|
| User Apps | User Apps | User Apps | User Apps |
| **USER** | | | |
| **KERNEL** | | | |
| VFS/File System | VFS/File System | VFS/File System | |
| Block Layer | Block Layer | Block Layer | **?** |
| Request Queue | Request Queue | | |
| SCSI XLAT | SCSI XLAT | | |
| SAS Driver | SAS Driver | NVMe Driver | |

Optimized Stack (SAS SSD) / Minimized Stack (NVMe SSD)

HDD (~10ms)  SAS SSD (~150μs)  NVMe SSD (<100μs)  PM (20~100ns)

NECSST
Next-generation Embedded / Computer System Software Technology

USER

KERNEL

User Apps

User Apps

User Apps

User Apps

VFS/File System

VFS/File System

VFS/File System

Since the medium is very fast,
the overhead of the **critical path** should be drastically reduced.

1. Multi-versioning mechanism

2. Long I/O stack

3. Management overhead

We need to rethink the critical path components that cause performance degradation.

SAS Driver

SAS Driver

NVMe Driver

HDD (~10ms)

SAS SSD (~150μs)

NVMe SSD (<100μs)

intel OPTANE DC
PERSISTENT MEMORY

PM (20~100ns)

SCHOOL OF ELECTRICAL AND
COMPUTER ENGINEERING

CISSR
Center for Innovative System Software Research

**NECSST**
**Next-generation Embedded / Computer System Software Technology**

- **PM-dedicated file system and tiered PM file system**

Next-generation Embedded / Computer System Software Technology

- **Designed to reap PM performance**



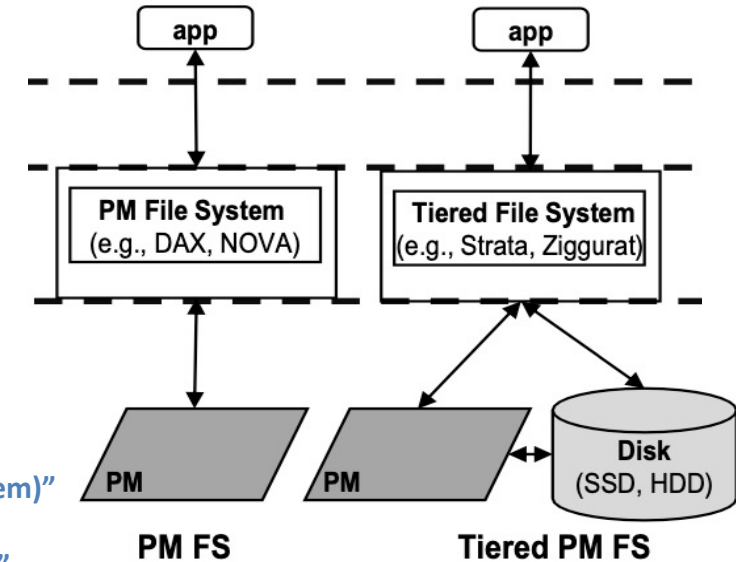| | |
|---|---|
| SOSP 2009 | "BPFS (Better I/O Through Byte-Addressable, Persistent Memory)" |
| SC 2011 | "SCMFS (SCMFS: A File System for Storage Class Memory)" |
| EuroSys 2014 | "PMFS (System Software for Persistent Memory)" |
| EuroSys 2014 | "Aerie (Aerie: Flexible File-System Interfaces to Storage-Class Memory)" |
| EuroSys 2016 | "HiNFS (A High Performance File System for Non-Volatile Main Memory)" |
| SOSP 2017 | "NOVA (NOVA-Fortis: A Fault-Tolerant Non-Volatile Main Memory File System)" |
| SOSP 2017 | "Strata (Strata: A Cross Media File System)" |
| HotStorage 2019 | "EvFS (EvFS: User-level, Event-driven File System for Non- volatile Memory)" |
| FAST 2019 | "Orion (Orion: A Distributed File System for Non-Volatile Main Memory and RDMA-Capable Networks)" |
| FAST 2019 | "Ziggurat (Ziggurat: A Tiered File System for Non- Volatile Main Memories and Disks)" |
| SOSP 2019 | "ZoFS (Performance and Protection in the ZoFS User-space NVM File System)" |
| SOSP 2019 | "SplitFS (SplitFS: Reducing Software Overhead in File Systems for Persistent Memory)" |
| FAST 2021 | "KucoFS (Scalable Persistent Memory File System with Kernel-Userspace Collaboration)" |

-------------------------------------------------------------------------------------------------------------------------------

| | |
|---|---|
| Linux kernel | "DAX (Ext4-DAX, XFS-DAX)" |

- **PM only**
  - PM as end destination media
  - Replace traditional storage?
    - Exception: Strata and Ziggurat

- **Lengthy process to maturity**
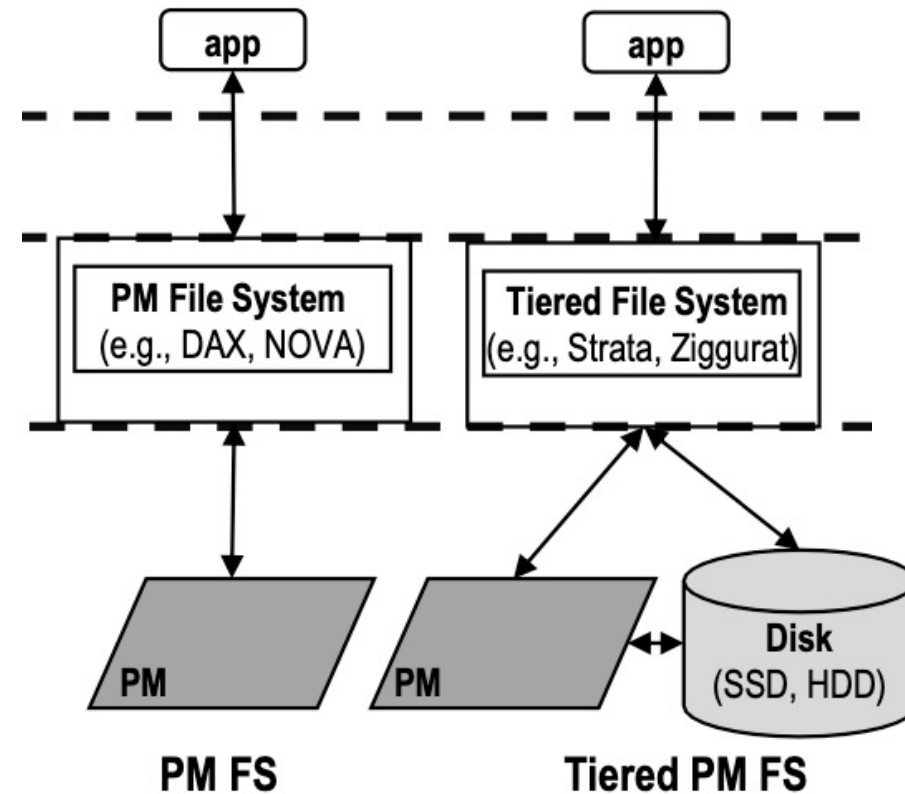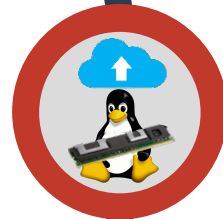  - E.g., Ext4…still in progress
  - Wisdom with age

Next-generation Embedded / Computer System Software Technology

# TABLE OF CONTENTS

NECSST

Next-generation Embedded / Computer System Software Technology

- **PM-based cache-like layer**



| Application |
| --- |

**Virtual File System**

**FR**
Intel OPTANE DC

**Other I/O Stacks**

**File System**
JFS EXT3 NILFS2
EXT4 BTRFS
FAT XFS F2FS NOVA
EXT2 Ext4-DAX
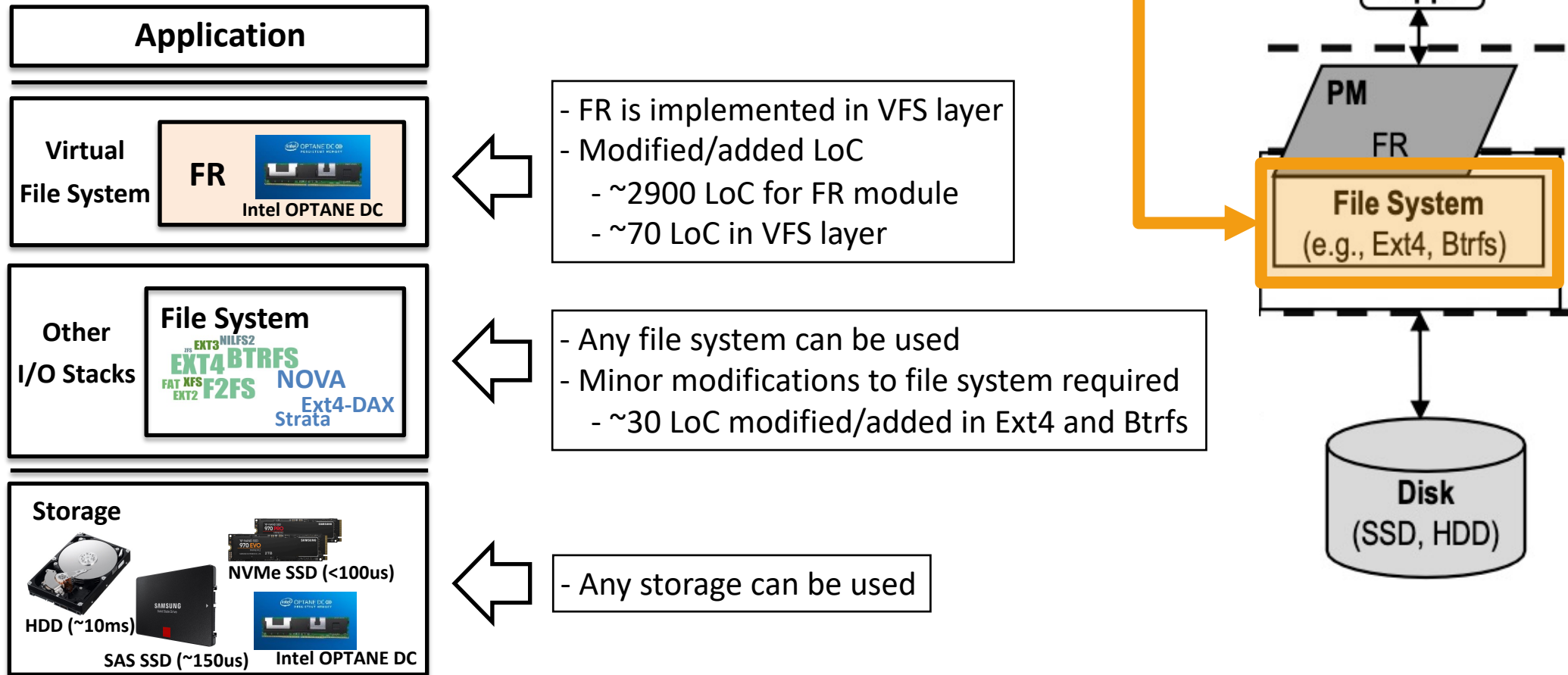Strata

**Storage**
HDD (~10ms)
SAS SSD (~150us)
NVMe SSD (<100us)
Intel OPTANE DC

- FR not only acts as a **cache**, but also as a **storage** using its persistent properties

- **Keep legacy file system and storage media "as-is"**



| Application |

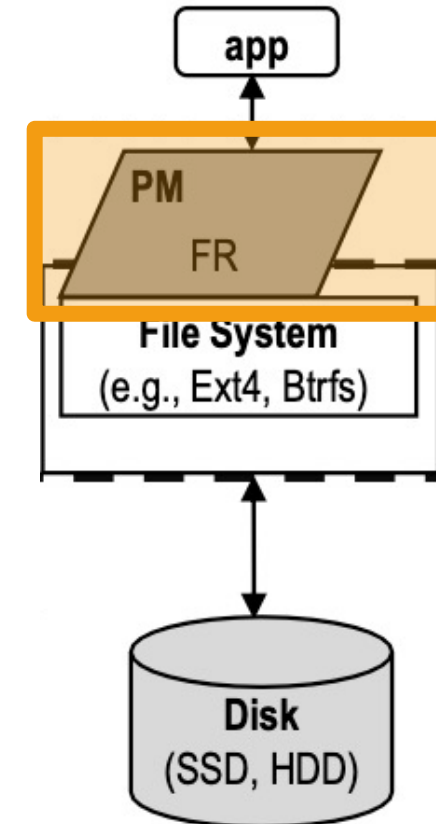**Virtual File System** — FR (Intel OPTANE DC)

- FR is implemented in VFS layer
- Modified/added LoC
    - ~2900 LoC for FR module
    - ~70 LoC in VFS layer

**Other I/O Stacks** — File System (ZFS, EXT3, NILFS2, EXT4, BTRFS, FAT, XFS, EXT2, F2FS, NOVA, Ext4-DAX, Strata)

- Any file system can be used
- Minor modifications to file system required
    - ~30 LoC modified/added in Ext4 and Btrfs

**Storage** — HDD (~10ms), SAS SSD (~150us), NVMe SSD (<100us), Intel OPTANE DC

- Any storage can be used

app

PM

FR

File System (e.g., Ext4, Btrfs)

Disk (SSD, HDD)

## PM performance

- Lightweight static management

**\* Average latency for managing cache for various indexing and management policies**

| Our static indexing | 67ns |
|---|---|
| Hashing indexing | 75ns |

VS.

| Activity | | Radix-tree + LRU | Hash + LRU |
|---|---|---|---|
| Radix-tree / Hash | Hashing | - | 78ns |
| | Search | 35ns | 162ns |
| | Insert* | $19\mu s$ | $19\mu s$ |
| | Delete | 190ns | 43ns |
| LRU | Touch | 161ns | 153ns |
| | Add | 73ns | 65ns |
| | Remove | 88ns | 82ns |

\* Insert includes mechanism to find empty blocks

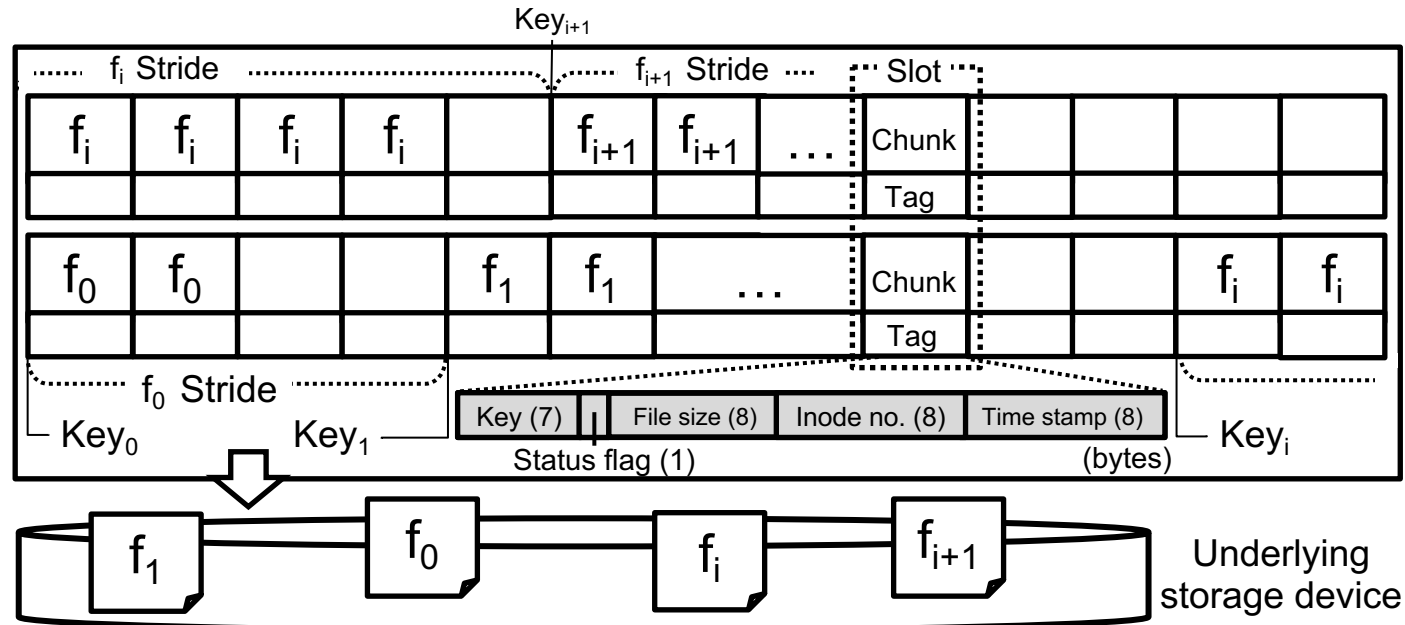## Ensure durability/consistency

- Static protocol naturally fulfills this

app

PM
FR

File System
(e.g., Ext4, Btrfs)
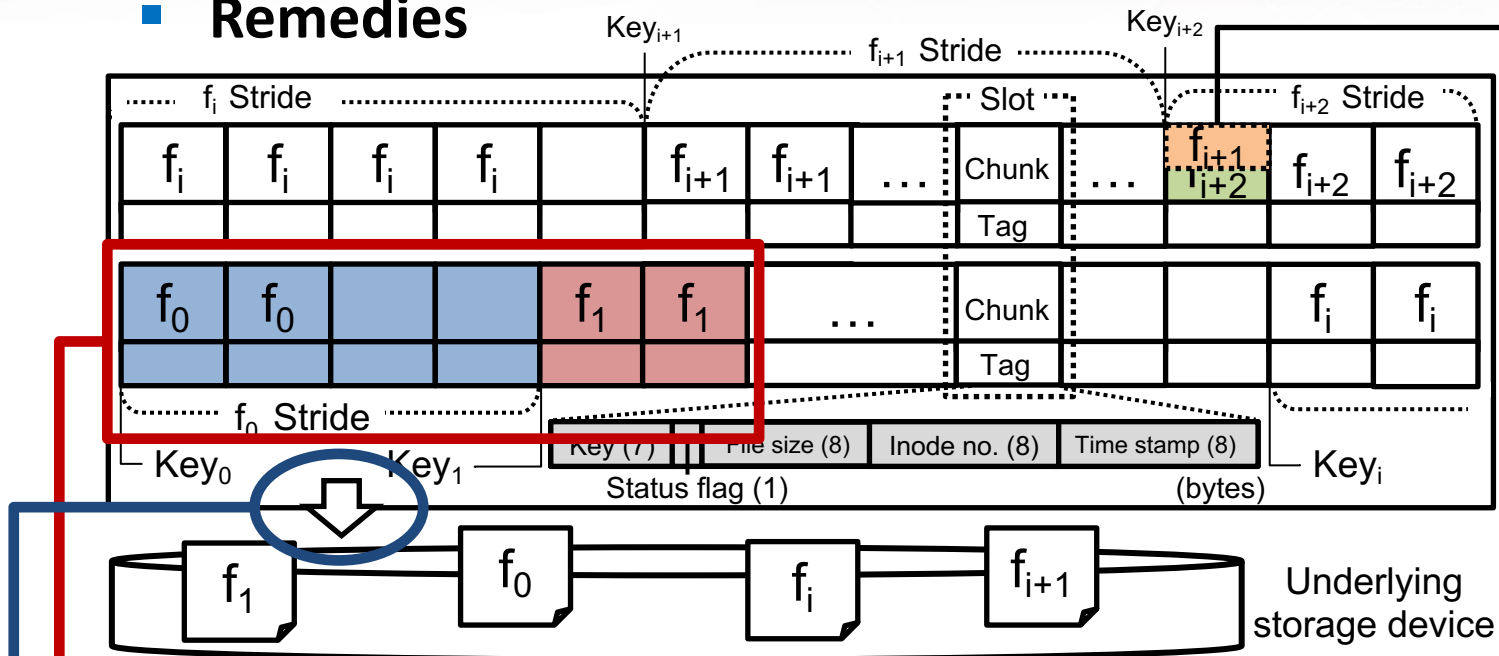
Disk
(SSD, HDD)

- **Internal components of FR**



- Chunk: Actual data is stored
- Tag: Some file's information and the status of chunk are stored
  - Key used for indexing: **key** mod Floor(N/2)
  - Bits in Status flag: V (Valid), N (New)

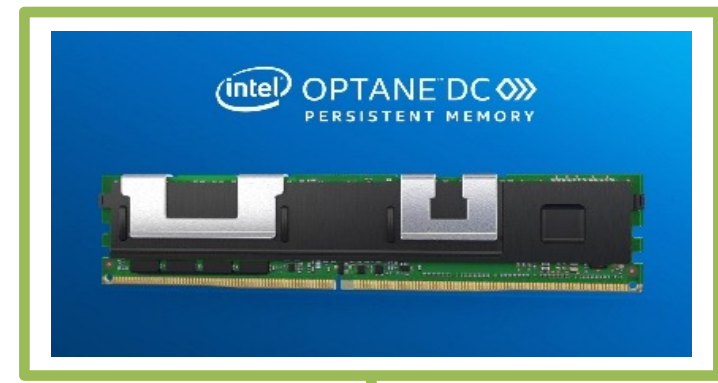Next-generation Embedded / Computer System Software Technology

- **Layout of FR**



- Static placement/replacement scheme in FR
  - Every files have destined location within FR with no PM allocator
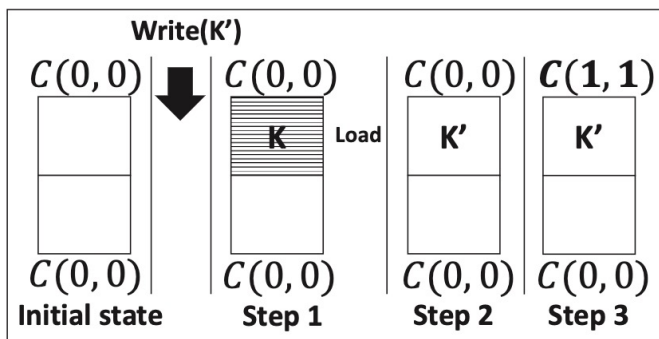    - → Can result in higher miss rate and collision

Next-generation Embedded / Computer System Software Technology

■ **Remedies**



*$Key_{i+1}$ $f_{i+1}$ Stride $Key_{i+2}$*

| $f_i$ Stride | | | | | | Slot | | $f_{i+2}$ Stride |
|---|---|---|---|---|---|---|---|---|

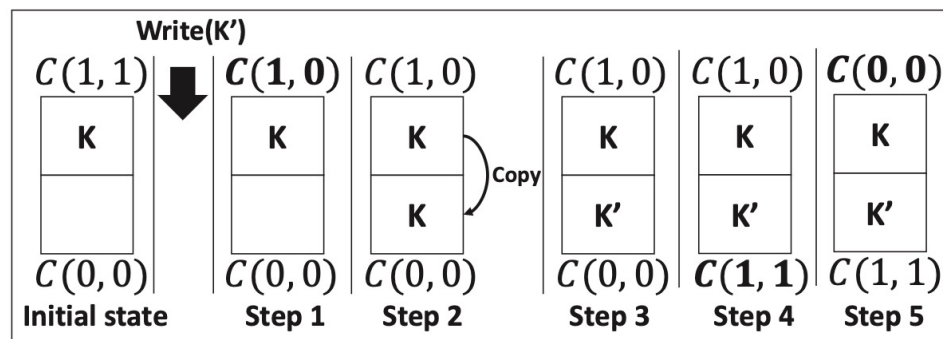**\*Collision** occurs when data from files invades each other's chunk beyond the stride value

* Sufficient large FR (>> working set)

* Stride: To eliminate invasion in chunk as much as possible

  − Files are positioned apart from each other by stride

* Periodic Flush: To reduce penalty (for clean chunk, there is no penalty for collision)

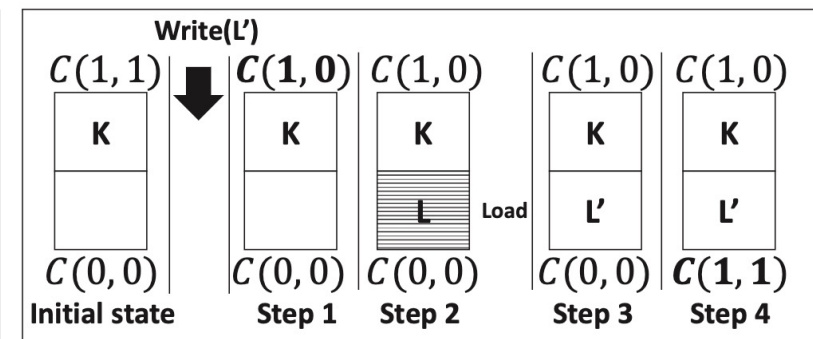  − Data is written to chunk, is flushed in the background

- **Steps taken based on initial condition of slots when write is requested**



(a) Case 1

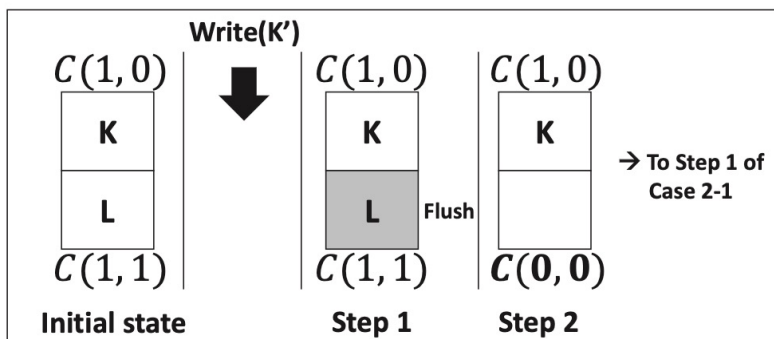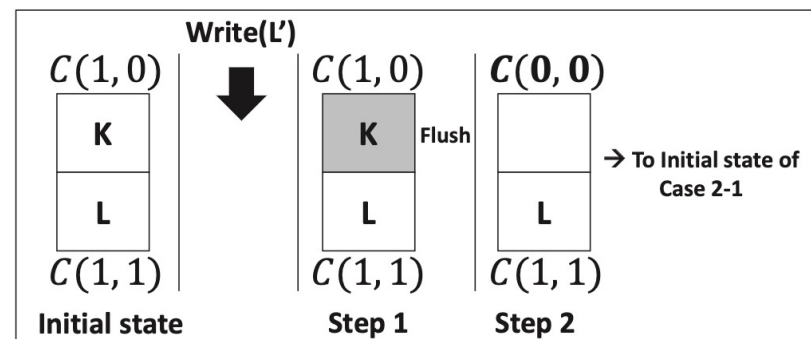(b) Case 2-1

(c) Case 2-2

(d) Case 3-1

(e) Case 3-2

(f) Case 3-3

# TABLE OF CONTENTS

Next-generation Embedded / Computer System Software Technology

## System configuration

| | Description |
|---|---|
| CPU | Intel(R) Xeon(R) Gold 6242 CPU @ 2.80GHz (16 cores) |
| DRAM | Samsung 32GB 2666MHz DDR4 RDIMM×4 (128GB) |
| PM | Intel Optane DC Persistent Memory (128GB) |
| Storage | Samsung V-NAND SSD 860 EVO (1TB) |
| OS | Linux Ubuntu 18.04.3 LTS (64bit) kernel v4.18 |



## Description of experimental comparison

| Notation | Description | Configuration | | |
|---|---|---|---|---|
| | | **PM** | **DRAM** | **Backing storage** |
| FR-X | FR applied to Ext4 (X is period value, e.g., 10ms) | 128GB | | 1TB (SSD) |
| Ext4 | Traditional block-based file system | | 128GB | 1TB (SSD) |
| DM-WC | DM-Writecache applied to Ext4 | 128GB | 128GB | 1TB (SSD) |
| DAX | PM-aware file system developed based on Ext4 | | | 128GB (PM) |
| NOVA | PM-aware file system | | | 128GB (PM) |

**NECSST**
*Next-generation Embedded / Computer System Software Technology*

- **Benchmarks**

| Filebench | R:W | Mean file size | # of files | Key-value store | Data set size | R:W |
|---|---|---|---|---|---|---|
| Fileserver | 1:2 | 128KB | 200K | YCSB-A,-F | 4GB | 1:1 |
| Varmail | 1:1 | 32KB | 800K | YCSB-B,-D,-E | 4GB | 19:1 |
| OLTP | 1:1 | 1.5GB | 20 | YCSB-C | 4GB | 1:0 |

- Filebench
  - Fileserver: write-intensive workload *without* fsync() calls
  - Varmail and OLTP: have considerable number of fsync() calls
- YCSB (record selection for -D is Latest, while all others are Zipfian)
  - Application: RocksDB
  - -A, -F: write-intensive workloads
  - -C: read-only workload
  - -B, -D, -E: read-intensive workloads

**NECSST**
Next-generation Embedded / Computer System Software Technology
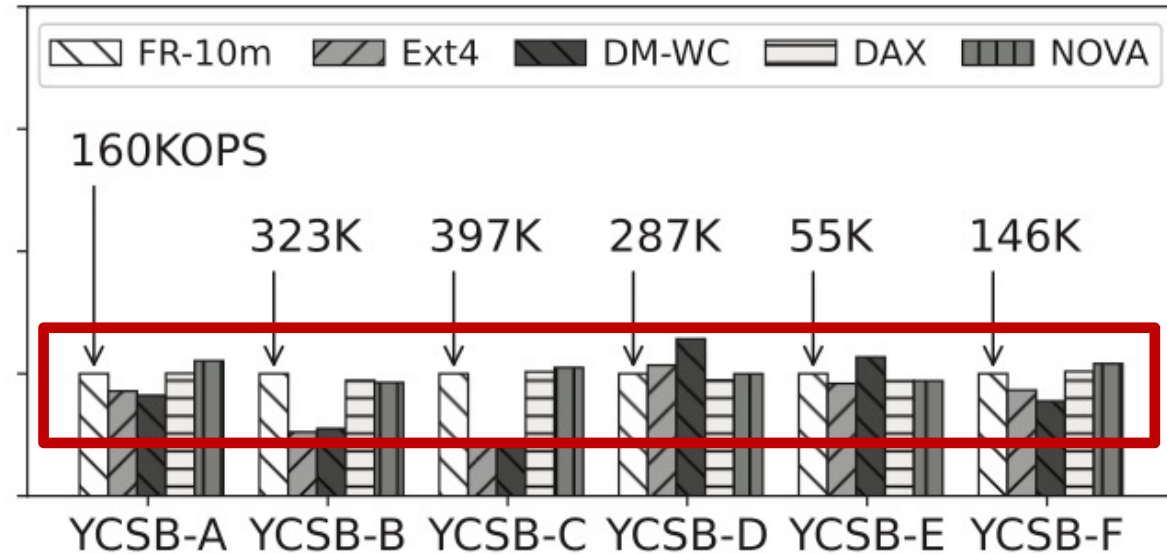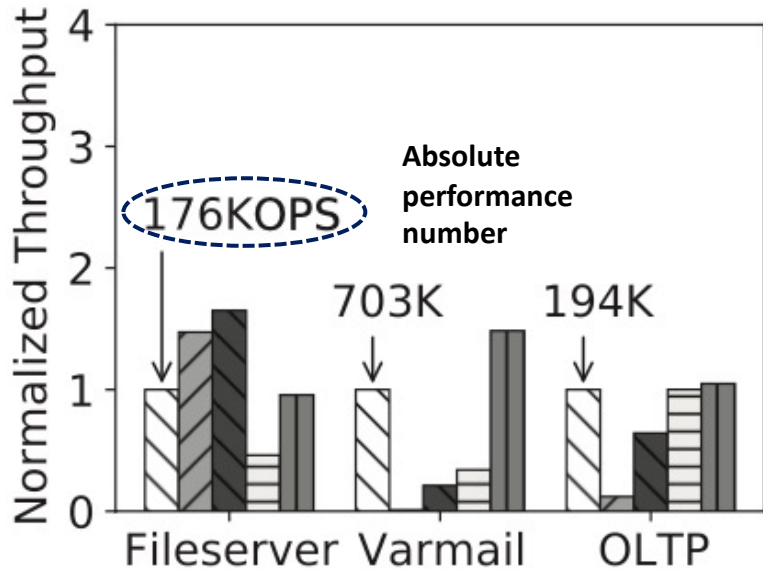
- **Overall performance**



**Observations**

- [Filebench] For Varmail and OLTP, FR is roughly 94x and 8x better than Ext4 and roughly 4.5x and 1.5x better than DM-WC
     FR performance is slightly lower than NOVA, while DAX suffers for Varmail
- [YCSB] FR, NOVA, and DAX show similar performance
     Ext4 and DM-WC perform worst for YCSB-A, -B, and -C

**NECSST**
Next-generation Embedded / Computer System Software Technology
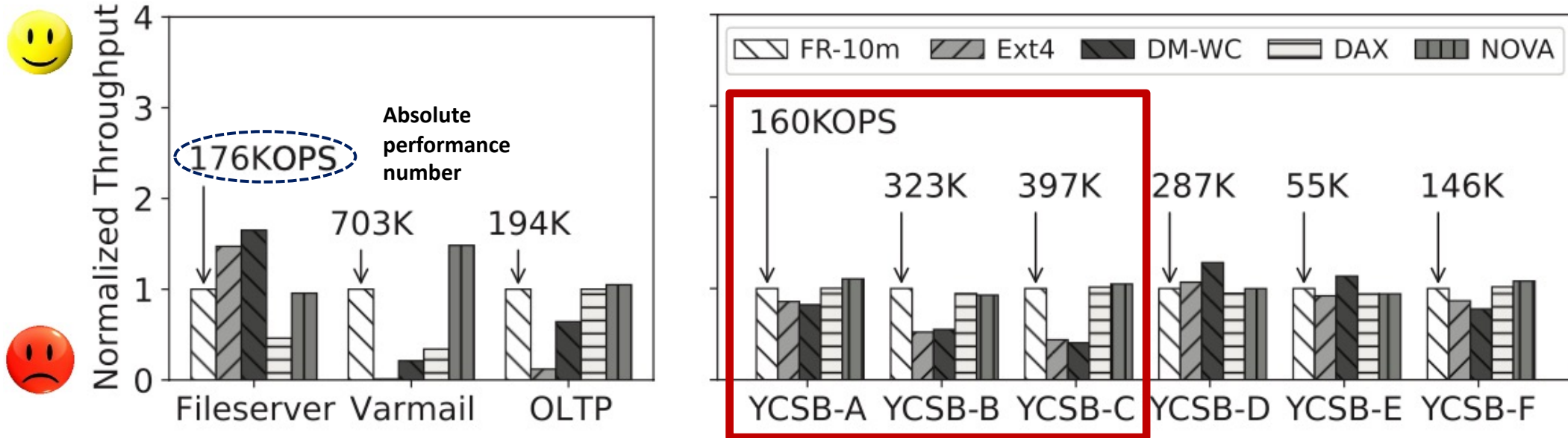
- **Overall performance**



**Observations**

- [Filebench] For Varmail and OLTP, FR is roughly 94x and 8x better than Ext4 and roughly 4.5x and 1.5x better than DM-WC

    FR performance is slightly lower than NOVA, while DAX suffers for Varmail

- [YCSB] FR, NOVA, and DAX show similar performance

    Ext4 and DM-WC perform worst for YCSB-A, -B, and -C

**NECSST**

Next-generation Embedded / Computer System Software Technology
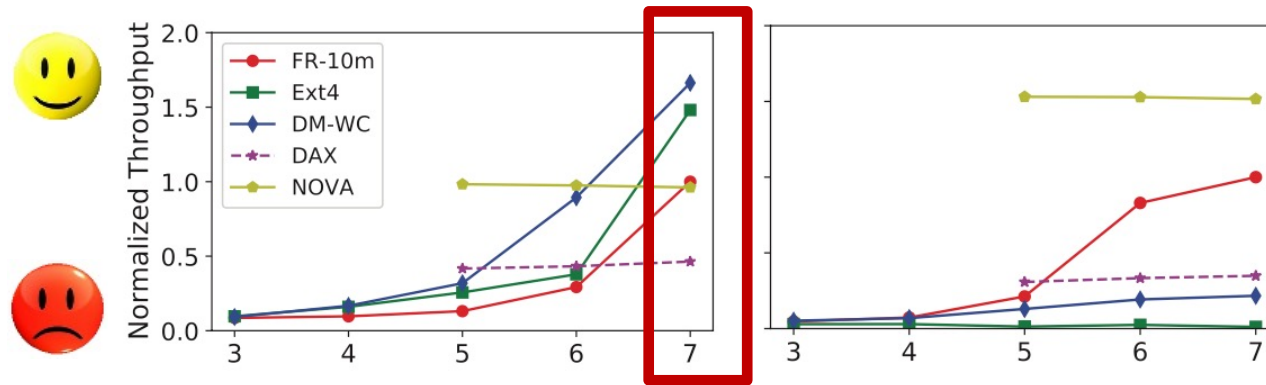
## Overall performance



**Observations**

- [Filebench] For Varmail and OLTP, FR is roughly 94x and 8x better than Ext4 and roughly 4.5x and 1.5x better than DM-WC

      FR performance is slightly lower than NOVA, while DAX suffers for Varmail

- [YCSB] FR, NOVA, and DAX show similar performance

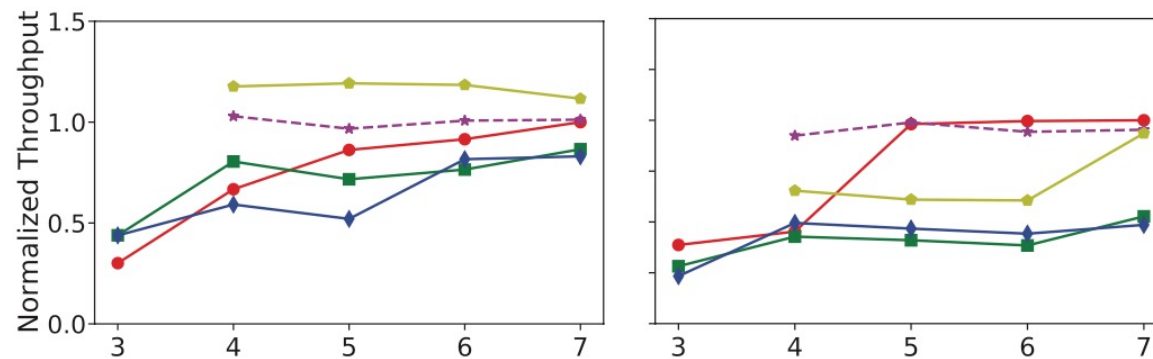      Ext4 and DM-WC perform worst for YCSB-A, -B, and -C

**Next-generation Embedded / Computer System Software Technology**

- **Performance results for PM size of $2^x$GB, where *x* is value of points in *x*-axis**
  - Normalized to the performance of FR when *x* = 7 (128GB)
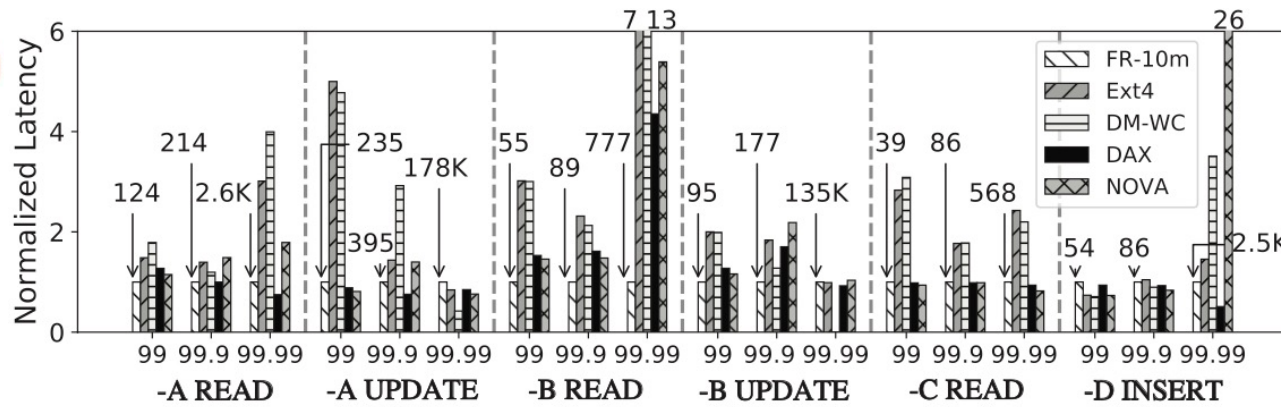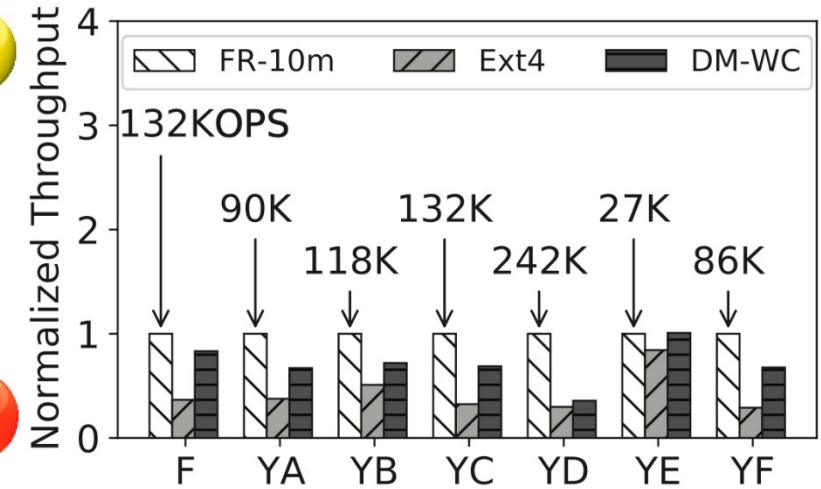


(a) Filebench
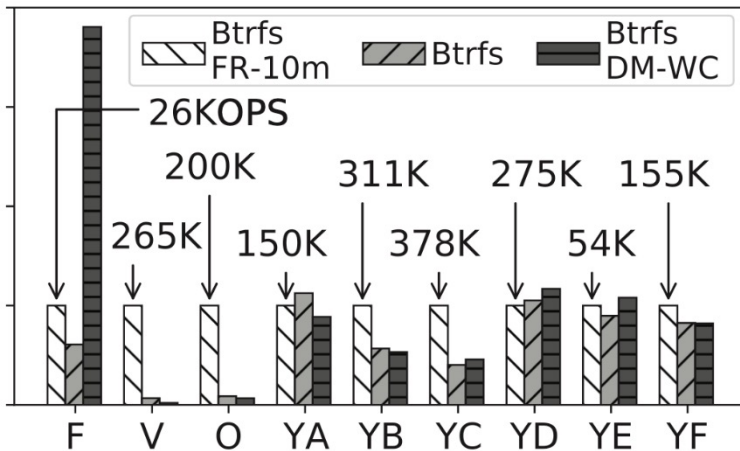
(b) Varmail

(c) YCSB-A

(d) YCSB-B

Next-generation Embedded / Computer System Software Technology

- **Tail latency**



- **Compensating for extra PM**



- **FR applied to Btrfs**



- **With different storage devices**



(a) NVMe SSD    (b) HDD

**Next-generation Embedded / Computer System Software Technology**

Read    Write

File ID and Offset

Times (sec)

**Working set of standard workload (Fileserver)**

■ **Limitations of standard workloads**

- Standard workloads do not capture the dynamics of real-world workloads
  - In terms of pattern
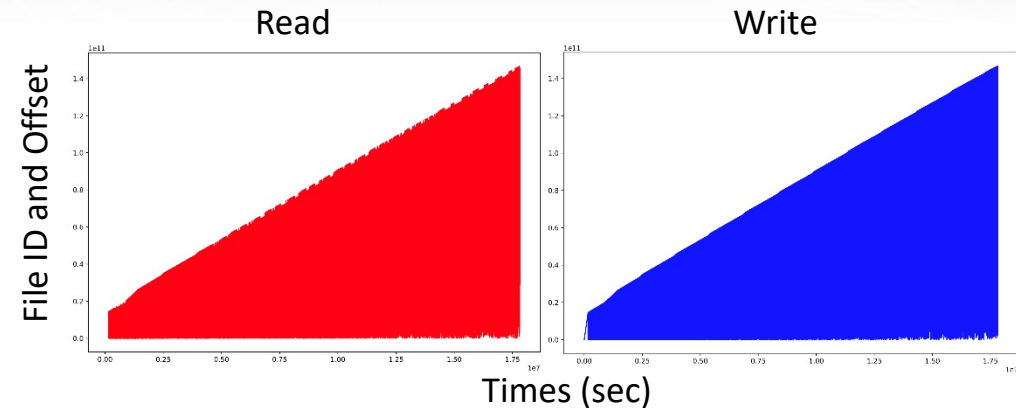    - Standard workloads: Working set does not change with time
    - Real-world workloads: Working set grows and shrinks as time evolves
  - In terms of operation generation
    - Standard workloads: No change in the access intensities of working set over time
    - Real-world workloads: Access intensities are also vary with time

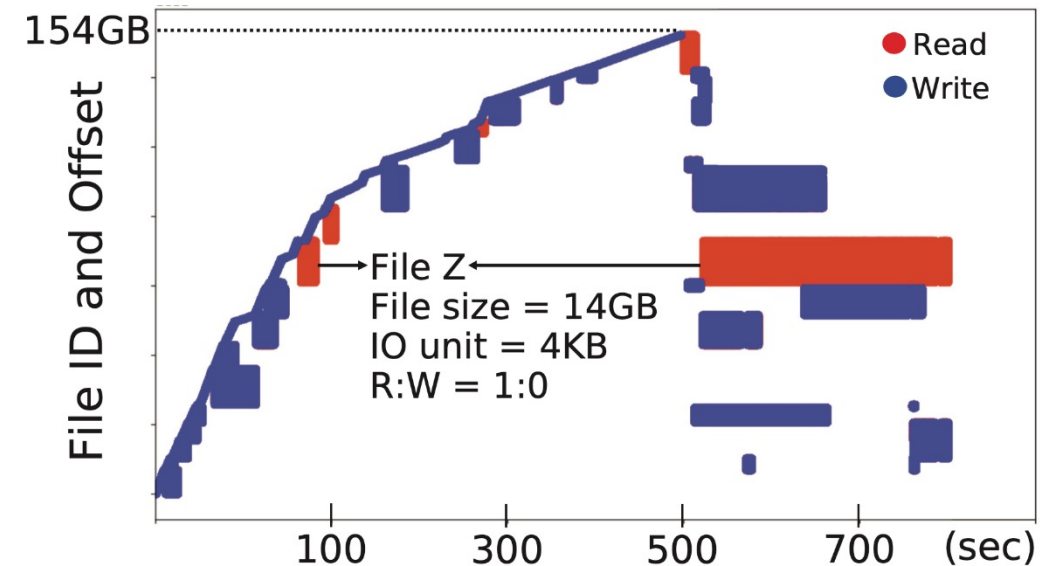→ Need for **dynamic workload** that is more representative of real-world workloads

**NECSST**
Next-generation Embedded / Computer System Software Technology

- **We devise synthetic workloads using I/O testing tool FIO**

| | Description |
|---|---|
| File size | 1~14GB [whole numbers only: 1 to 2 files of each] |
| Distribution | Pareto 0.1/0.5 [2/2], Zipf 0.2/1.2 [5/1], N [6], R [5] |
| IO unit | 4KB [14], 8KB [7] |
| Read:Write | 1:0 [4], 19:1 [9], 1:1 [8] |
| Fsync | 0~5% [whole numbers only: 2 to 5 files of each] |
| Intensity | Request interval: 1$\mu$s~1ms (randomly distributed) |

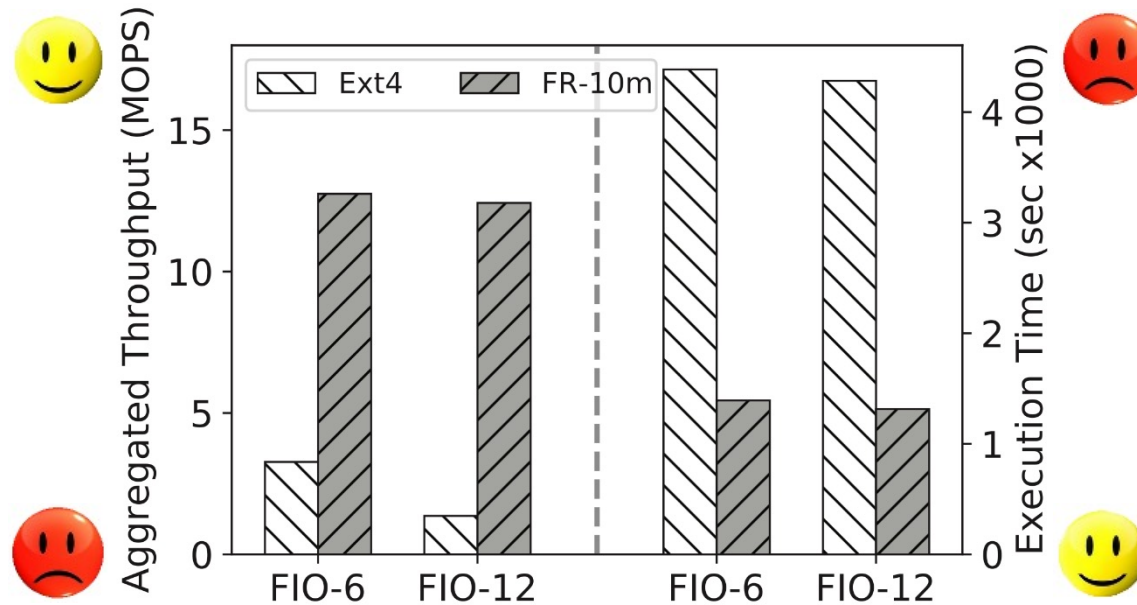(a) Characteristics of 21 files used to generate synthetic workload (FIO)



(b) Working set (FIO-12)

**Configuration: Total IO size is 575GB**
- FIO-6: 6 files, 50GB working set
- FIO-12: 12 files, 100GB working set

Next-generation Embedded / Computer System Software Technology

- **Performance results**



**Observations**
- FR provides more than 9x higher aggregate throughput and ended over 3x faster than Ext4
- FR is providing immediately durable in-order semantics
- For NOVA and DAX, cannot run as dataset is larger than PM size

Next-generation Embedded / Computer System Software Technology

# TABLE OF CONTENTS

Next-generation Embedded / Computer System Software Technology

- **First Responder (FR)**
  - PM-based cache-like layer
  - Keep legacy file system and storage media "as-is"
  - PM performance
    - Respond quickly with in-order file system semantics
    - Hide traditional I/O stack overhead
  - Ensure durability/consistency
    - Protocol implemented with static management

Next-generation Embedded / Computer System Software Technology

hssong1987@unist.ac.kr

**NECSST**
Next-generation Embedded / Computer System Software Technology