# MadFS: Per-File Virtualization for Userspace Persistent Memory Filesystems
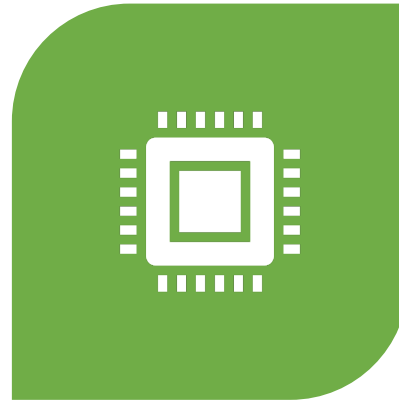
Shawn Zhong*, Chenhao Ye*, Guanzhou Hu, Suyan Qu

Andrea Arpaci-Dusseau, Remzi Arpaci-Dusseau, Michael Swift

WISCONSIN
UNIVERSITY OF WISCONSIN–MADISON

* Equal contribution

1

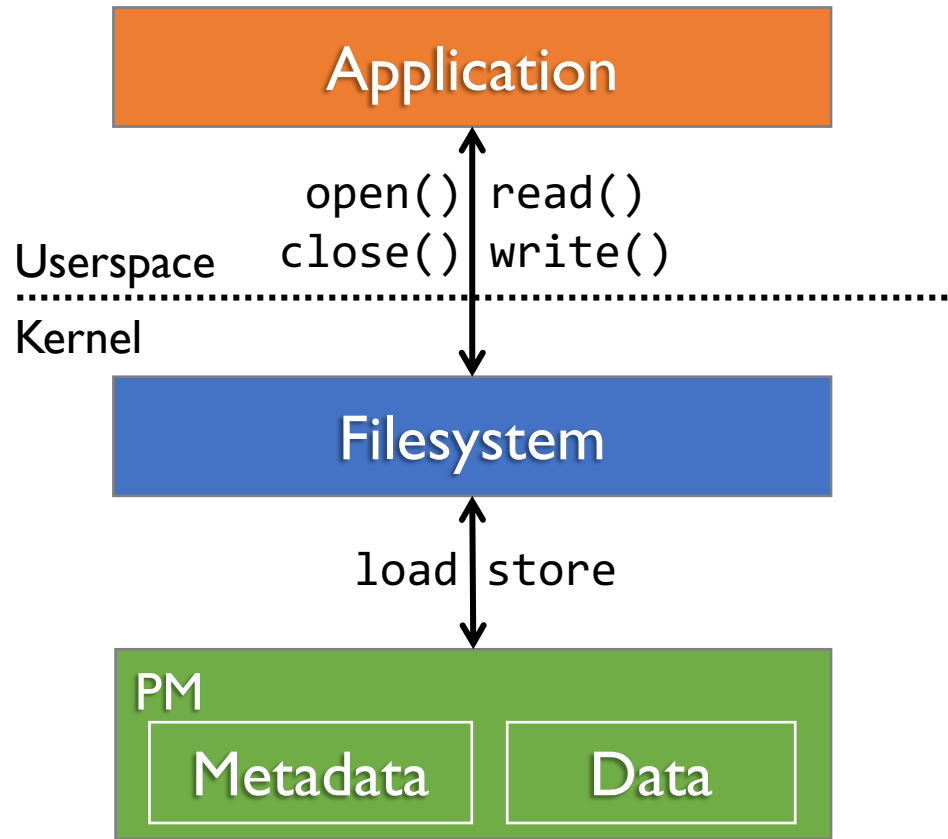# Background: Persistent Memory



*Fast*
sub-us latency

*Byte-Addressable*
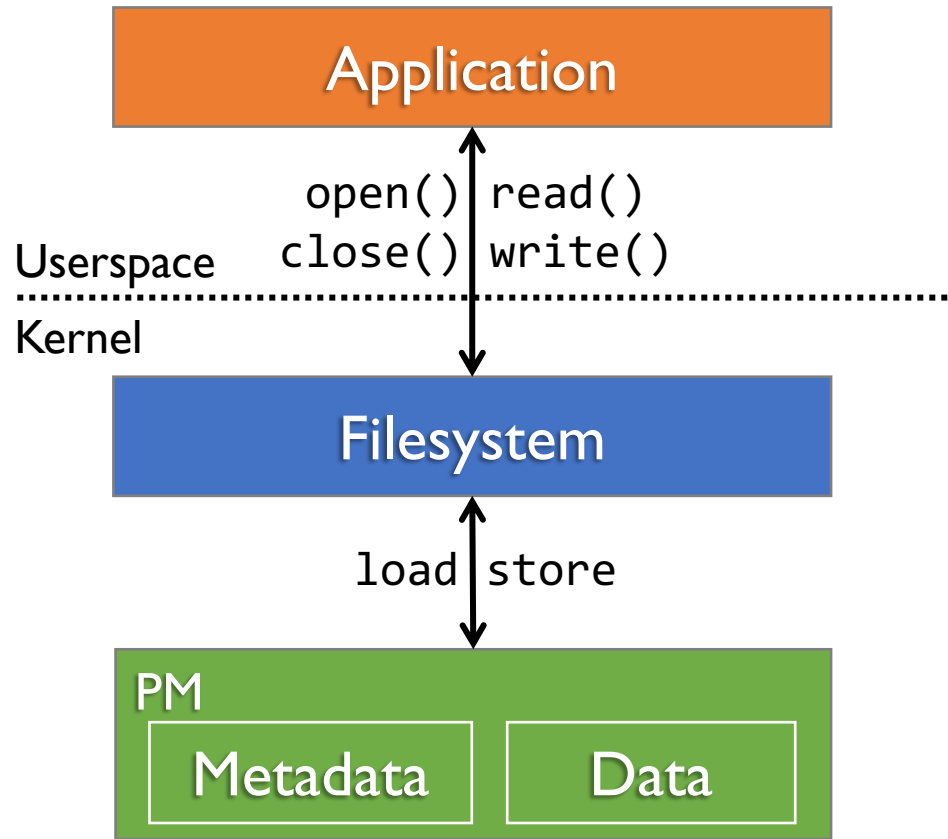accessed via CPU
instructions

*Non-Volatile*
retain data
without power

# Background: Kernel Filesystems for PM



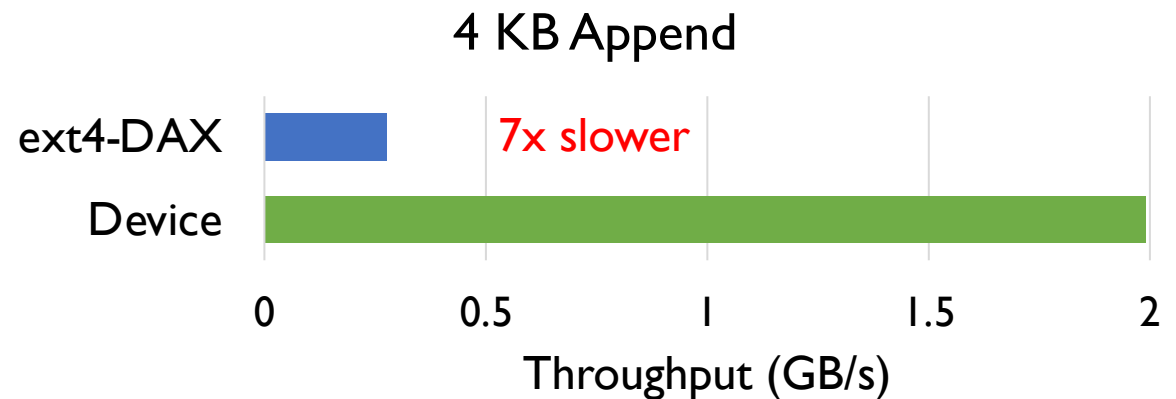Kernel FS manages both data & metadata

# Background: Kernel Filesystems for PM

Application

open() read()
close() write()

Userspace
Kernel

Filesystem
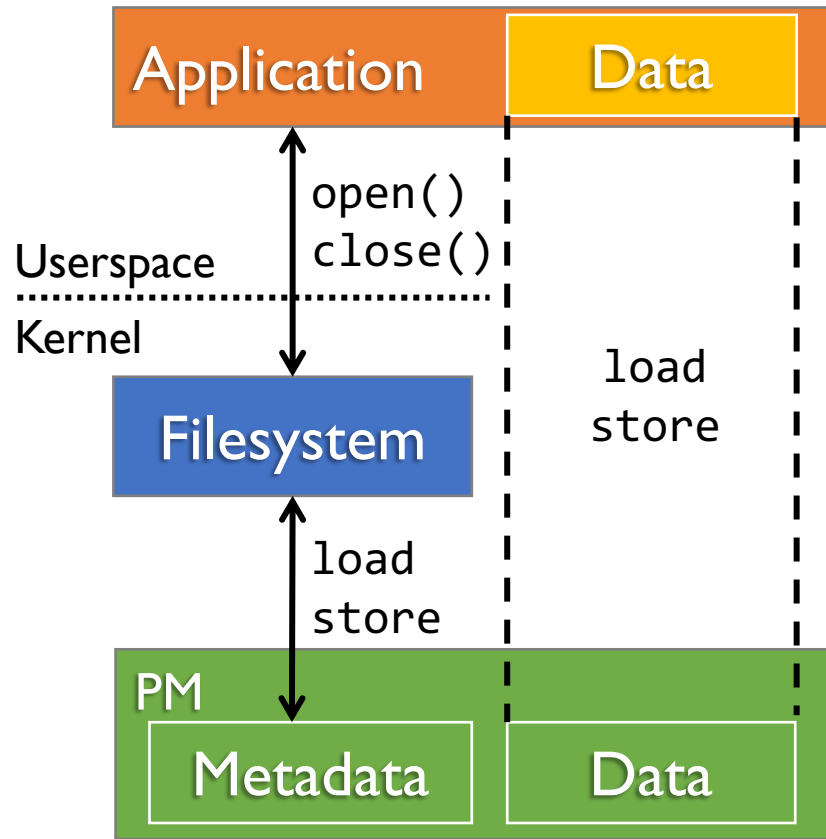
load store

PM

Metadata  Data

Kernel FS manages both data & metadata

Overhead for append in ext4-DAX

- System call

- VFS (e.g., inode locking)

- Metadata journaling in block granularity

### 4 KB Append

| | |
|---|---|
| ext4-DAX | 7x slower |
| Device | |

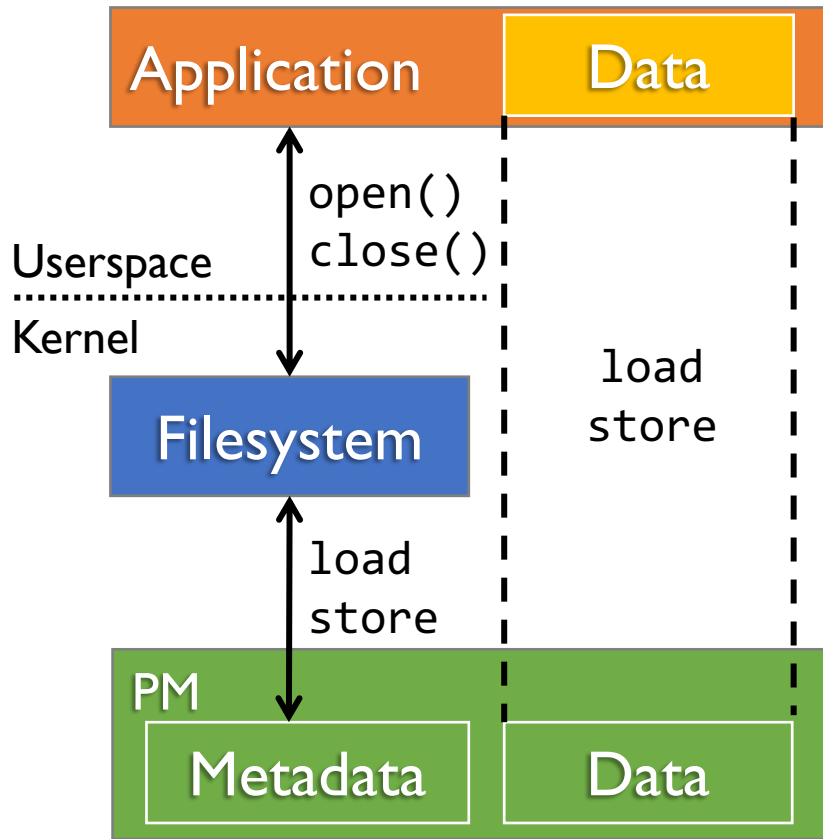Throughput (GB/s)

0    0.5    1    1.5    2

4

# Background: Userspace Filesystems for PM



Userspace FS bypass kernel for data ops

- Memory-map file data on open

- Handle read/write in userspace via load/store

# Background: Userspace Filesystems for PM



Userspace FS bypass kernel for data ops

- Memory-map file data on open

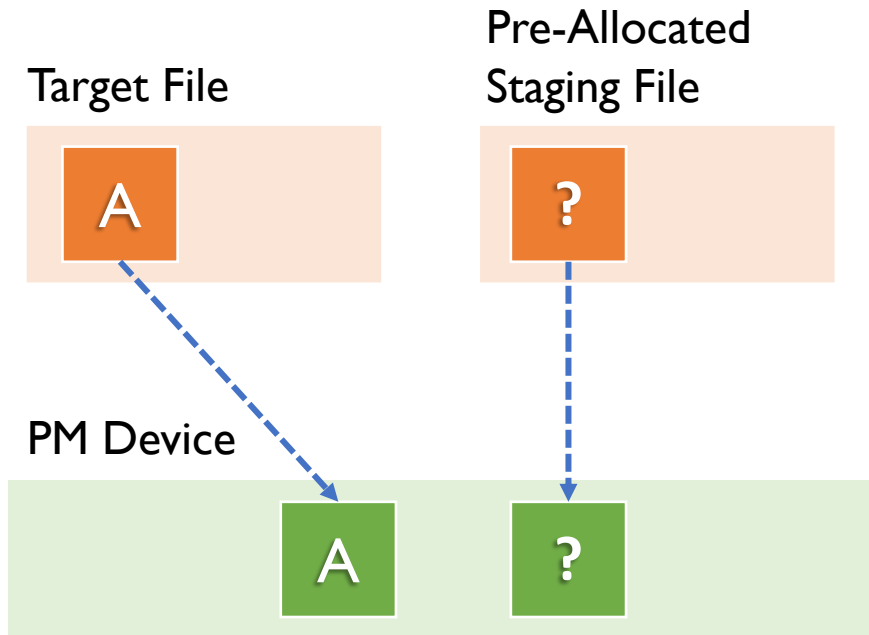- Handle read/write in userspace via load/store

Metadata still managed by kernel

*Issue*: Data ops coupled metadata updates

*Example*: Append + Fsync in SplitFS [SOSP '19]
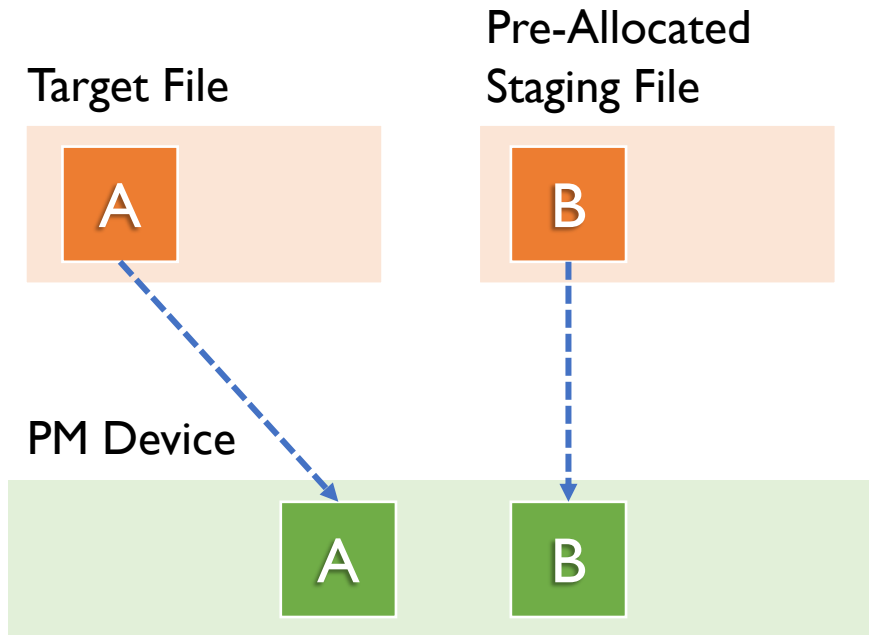
6

# Example: Append + Fsync in SplitFS
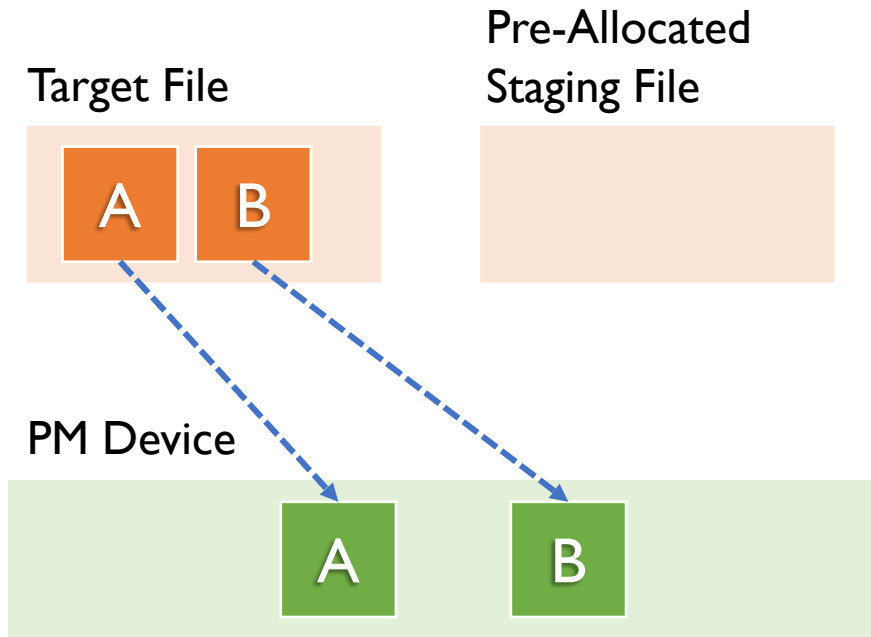
Append "B": Userspace Data Operation

Target File

Pre-Allocated
Staging File

PM Device

# Example: Append + Fsync in SplitFS



Append "B": Userspace Data Operation

- Write data to pre-allocated file

# Example: Append + Fsync in SplitFS

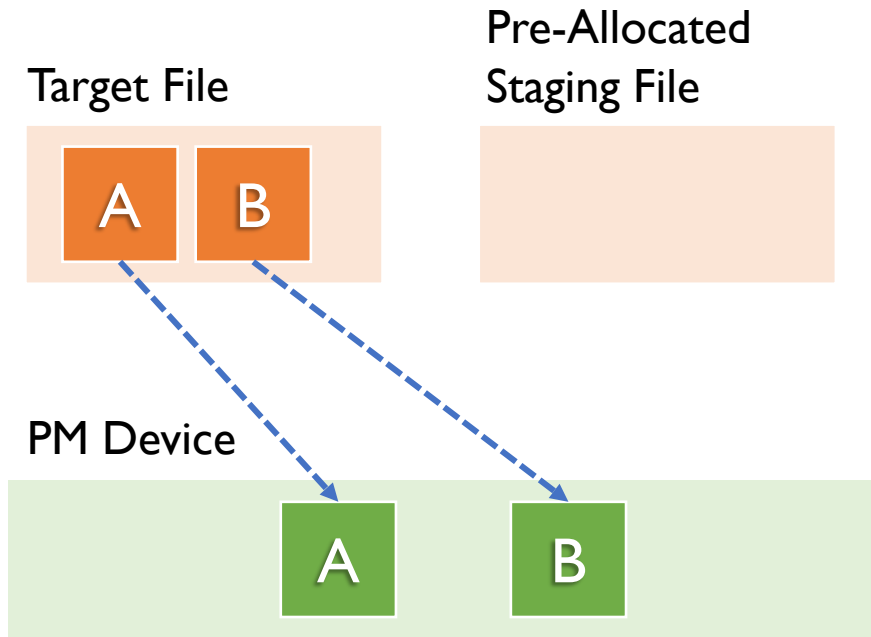Target File

Pre-Allocated
Staging File

A  B

PM Device

A  B

Append "B": Userspace Data Operation

• Write data to pre-allocated file

Fsync: Kernel Metadata Operation

• Remap data to target file for visibility

• Update block map in inode
  memory map / page table

# Example: Append + Fsync in SplitFS

Target File

Pre-Allocated Staging File
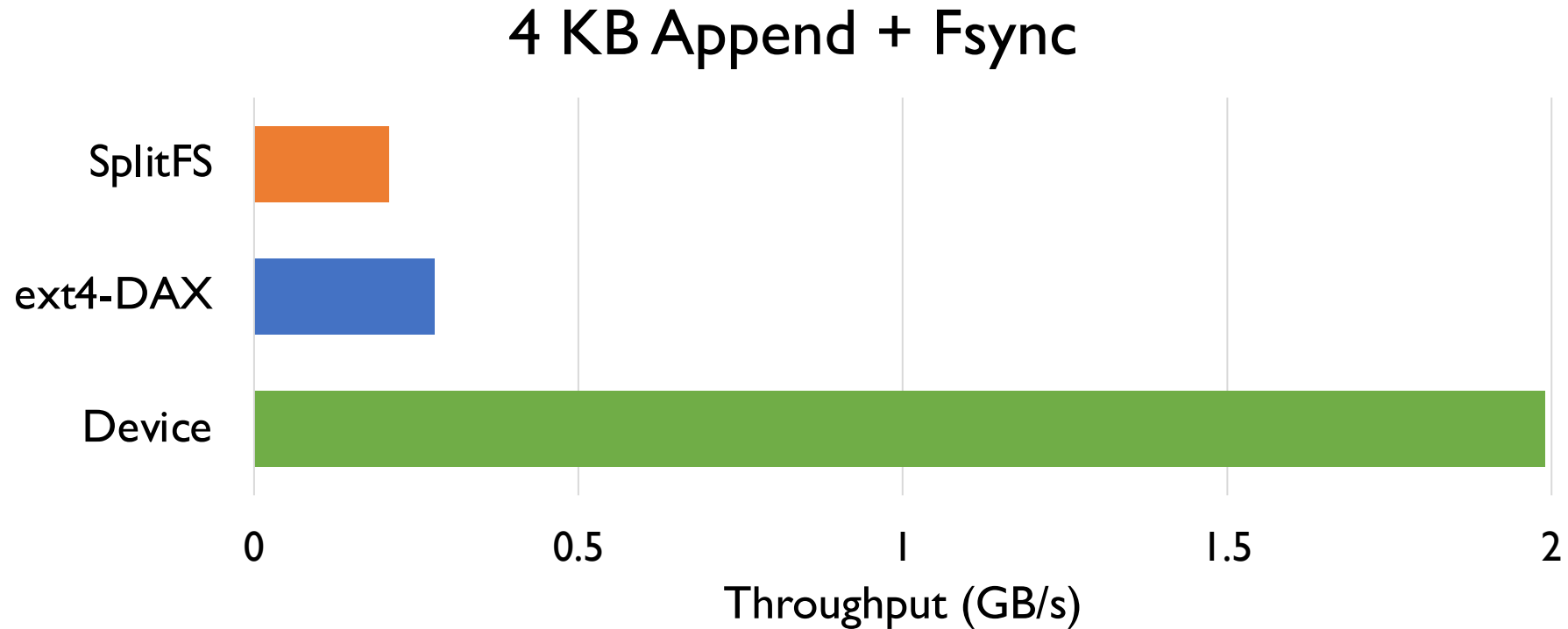
A  B

PM Device

A  B

Append "B": Userspace Data Operation

- Write data to pre-allocated file

Fsync: Kernel Metadata Operation

- Remap data to target file for visibility

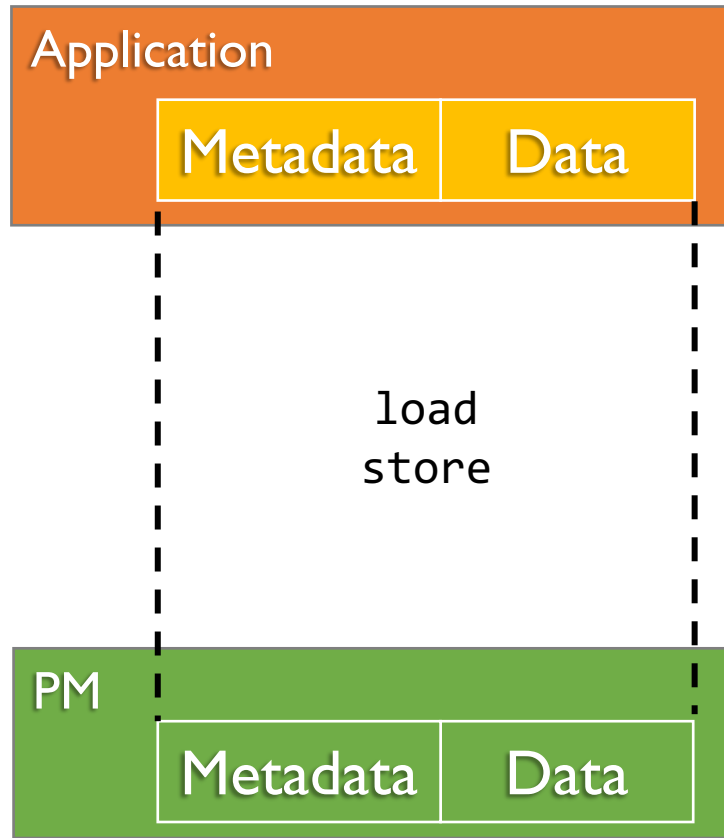- Update block map in inode: kernel I/O stack
  memory map / page table: TLB

Kernel metadata operation is expensive

10

# Example: Append + Fsync in SplitFS

## 4 KB Append + Fsync



*Result*: Worse performance compared to kernel FS 😢

# Background: Userspace Filesystems for PM

**Application**

| Metadata | Data |
|---|---|

load
store

**PM**

| Metadata | Data |
|---|---|

Expensive to modify kernel-managed metadata

Can we manage all metadata in userspace?

# Background: Userspace Filesystems for PM

Application

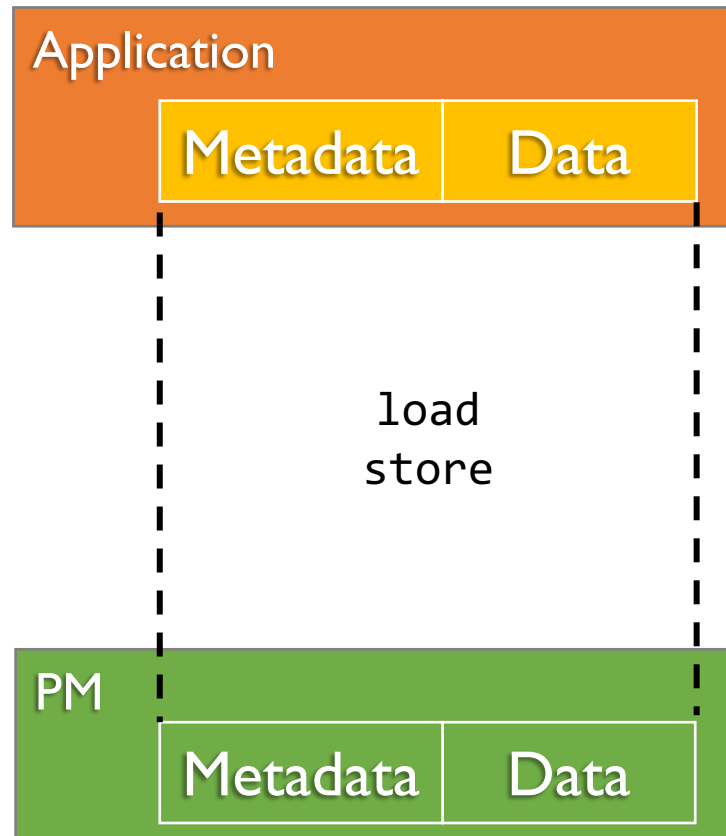| Metadata | Data |

load
store

PM

| Metadata | Data |

Expensive to modify kernel-managed metadata

Can we manage all metadata in userspace?

Unfortunately, no: applications are untrusted

_Example_: malicious user changes permission

# Background: Userspace Filesystems for PM



Expensive to modify kernel-managed metadata

Can we manage all metadata in userspace?

Unfortunately, no: applications are untrusted

*Example*: malicious user changes permission

What about only the metadata coupled with data operation

# 🔍 Observation:

## Some file metadata share

## the same protection domain as data



| Data | Size & Block Map | | Allocation & Owner |
| --- | --- | --- | --- |
| Userspace | | | Kernel |

💡 Insight:

Embed these metadata into data

for efficient userspace management
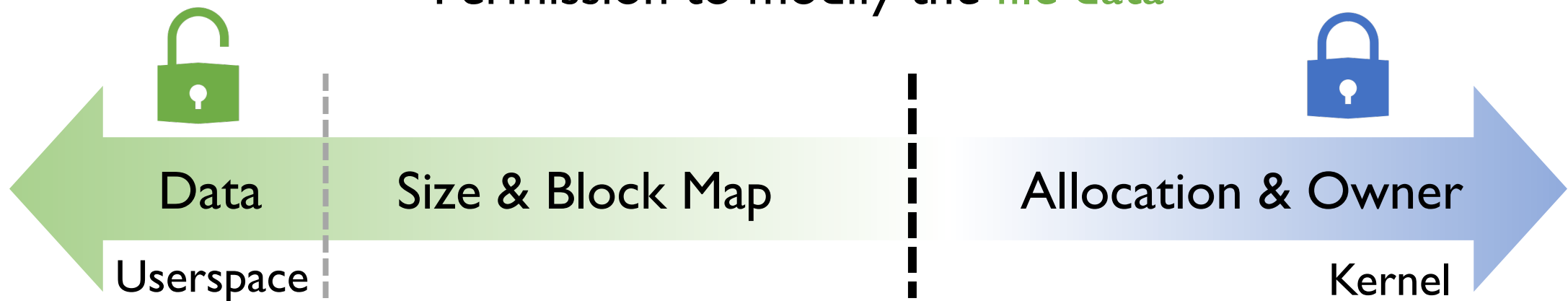
without sacrificing security

# Metadata Embedding

*Observation*: Some metadata share the same protection domain as data

*Example*: Block map

Permission to swap two block pointers within a file

≈

Permission to modify the file data



Data | Size & Block Map | Allocation & Owner

Userspace | | Kernel

17

# Metadata Embedding

*Insight*: Embed metadata coupled with data ops into file data

Efficient metadata operations

- No kernel I/O stack involvement

Equivalent security guarantees

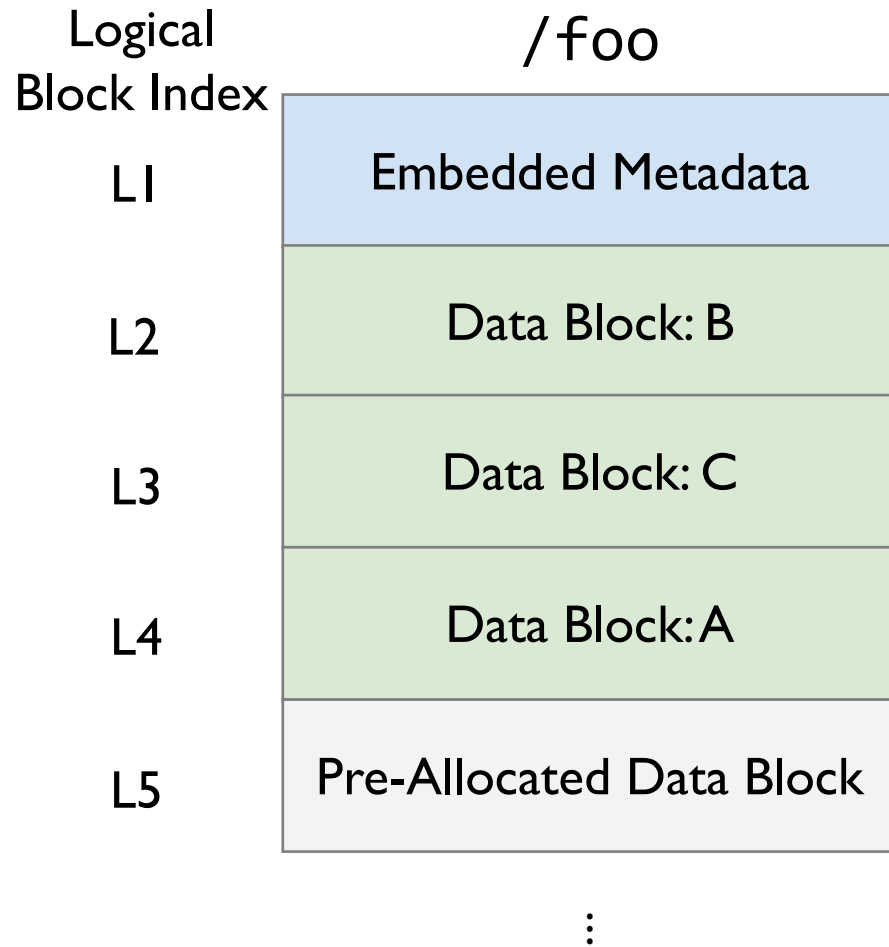- Require write permission to modify embedded metadata

# MadFS: Metadata Embedded Filesystem

Userspace library filesystem for PM

- Memory mapped I/O

- Data & most metadata ops in userspace
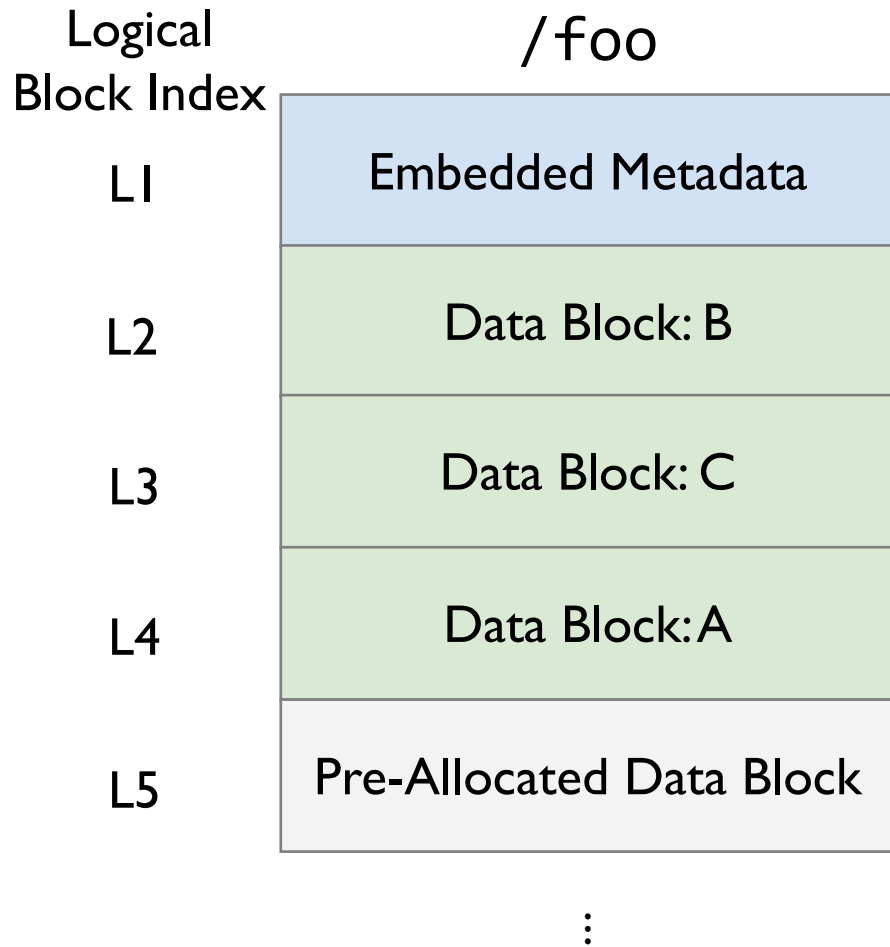
- Data crash consistency via copy-on-write

| POSIX Application | | | | |
|---|---|---|---|---|
| ↓ open | ↓ [p]read | ↓ [p]write | ↓ [f]stat | ↓ close |

MadFS Lib

. . .

| ↓ open | ↓ mmap | ↓ fallocate | ↓ fstat | ↓ close |
|---|---|---|---|---|

Userspace
Kernel

File on Unmodified Kernel FS (e.g., ext4-DAX)

# MadFS: Simplified Design

Logical Block Index

/foo

| | |
|---|---|
| L1 | Embedded Metadata |
| L2 | Data Block: B |
| L3 | Data Block: C |
| L4 | Data Block: A |
| L5 | Pre-Allocated Data Block |

⋮

Logical blocks: stored on the underlying FS

Virtual blocks: seen by the application

# MadFS: Simplified Design

Logical
Block Index

/foo

| | |
|---|---|
| L1 | Embedded Metadata |
| L2 | Data Block: B |
| L3 | Data Block: C |
| L4 | Data Block: A |
| L5 | Pre-Allocated Data Block |

⋮

Logical blocks: stored on the underlying FS

Virtual blocks: seen by the application

V1  V2  V3

| A | B | C |
|---|---|---|

Virtual

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

Logical

| M | B | C | A | ? | … |
|---|---|---|---|---|---|

L1  L2  L3  L4  L5

Block Map: {V1:L4, V2:L2, V3:L3}

Bitmap: 111100…

Virtual File Size: 12 KB

# MadFS: Simplified Design

Example: `pwrite(fd, buf, count=6KB, offset=10KB)`



V1    V2    V3    V4

| A | B | C | buf |

Virtual

Logical

| M | B | C | A | ? | ? | … |

L1   L2   L3   L4   L5   L6

Block Map: {V1:L4, V2:L2, V3:L3}

Bitmap: 111100…

Virtual File Size: 12 KB

# MadFS: Simplified Design

Example: `pwrite(fd,` `buf` `, count=6KB, offset=10KB)`

1. Allocate 2 logical blocks from the bitmap

| V1 | V2 | V3 |
|----|----|----|
| A  | B  | C  |

Virtual

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

Logical

| M | B | C | A | ? | ? | ⋯ |
|---|---|---|---|---|---|---|

L1  L2  L3  L4  L5  L6

Block Map: {V1:L4, V2:L2, V3:L3}

Bitmap: 111111⋯

Virtual File Size: 12 KB

# MadFS: Simplified Design

Example: `pwrite(fd,` buf `, count=6KB, offset=10KB)`

1. Allocate 2 logical blocks from the bitmap

2. Copy buffer and unaligned data

V1   V2   V3

| A | B | C |

Virtual

Logical

| M | B | C | A | C | buf | ... |

L1   L2   L3   L4   L5   L6

Block Map: {V1:L4, V2:L2, V3:L3}

Bitmap: 111111···

Virtual File Size: 12 KB

# MadFS: Simplified Design

Example: `pwrite(fd, buf, count=6KB, offset=10KB)`

1. Allocate 2 logical blocks from the bitmap

2. Copy buffer and unaligned data

3. Update block map and virtual size



Block Map: {V1:L4, V2:L2, V3:L5, V4:L6}

Bitmap: 111111⋯

Virtual File Size: 16 KB

# MadFS: Simplified Design

Example: `pwrite(fd,` `buf` `, count=6KB, offset=10KB)`

1. Allocate 2 logical blocks from the bitmap

2. Copy buffer and unaligned data

3. Update block map and virtual size

4. Deallocate old blocks

V1    V2    V3    V4

| A | B | C | buf |

Virtual

Logical

| M | B | C | A | C | buf | ... |

L1   L2   L3   L4   L5   L6

Block Map: {V1:L4, V2:L2, V3:L5, V4:L6}

Bitmap: 110111…

Virtual File Size: 16 KB

# MadFS: Simplified Design

Example: `pwrite(fd,` `buf` `, count=6KB, offset=10KB)`

1. Allocate 2 logical blocks from the bitmap

2. Copy buffer and unaligned data

3. Update block map and virtual size

4. Deallocate old blocks

Copy-on-write & append in userspace

V1  V2  V3  V4

| A | B | C | buf |

Virtual

Logical

| M | B | C | A | C | buf | ... |

L1  L2  L3  L4  L5  L6

Block Map: {V1:L4, V2:L2, V3:L5, V4:L6}

Bitmap: 110111…

Virtual File Size: 16 KB

# MadFS: Metadata Management

## Embedded Metadata

## Kernel-Managed Metadata

- Virtual-to-logical map

- Virtual file size

- Logical blocks bitmap

Allows efficient data ops without
expensive kernel involvement

# MadFS: Metadata Management

## Embedded Metadata

- Virtual-to-logical map

- Virtual file size

- Logical blocks bitmap

Allows efficient data ops without expensive kernel involvement

## Kernel-Managed Metadata

- Logical-to-physical map

- Logical file size

- Physical blocks bitmap

Updated on pre-allocation (infrequent)

- File permission

Provides coarse-grained allocation and protection

# MadFS: Per-File Virtualization

A userspace virtualization layer implements a complete set of file functionalities, including metadata management, crash consistency, and concurrency control, on a per-file basis



Metadata Management          Crash Consistency          Concurrency Control

# MadFS: Full Design (Details in Paper)

💥 Metadata Crash Consistency

- Log-structured metadata with 8-byte log entries

8-byte log entry

| LE | LE | LE |
|----|----|----|

Virtual 3-4 → Logical 5-6

# MadFS: Full Design (Details in Paper)

💥 Metadata Crash Consistency

- Log-structured metadata with 8-byte log entries

🔀 Lock-Free Optimistic Concurrency Control

- Commit log entry via compare-and-swap (CAS)

- Safe in presence of process crashes

- Better scalability with concurrent data ops

8-byte log entry

LE | LE | LE

Virtual 3-4 → Logical 5-6

CAS

CAS | commit?

LE | LE | LE |

# MadFS: Full Design (Details in Paper)

💥 Metadata Crash Consistency

- Log-structured metadata with 8-byte log entries

🔀 Lock-Free Optimistic Concurrency Control

- Commit log entry via compare-and-swap (CAS)

- Safe in presence of process crashes

- Better scalability with concurrent data ops

♻️ Non-Blocking Garbage Collection

- Read-Copy Update w/o tail latency impact

8-byte log entry

| LE | LE | LE |

Virtual 3-4 → Logical 5-6

CAS

CAS    commit?

| LE | LE | LE | |

head
❌ | LE | LE | LE | LE |
✅ | LE | LE |

# Evaluation

Questions:

- How does MadFS perform on microbenchmarks?

- How does MadFS perform on real-world applications?

Compare MadFS running on ext4-DAX with

- ext4-DAX, NOVA [FAST '16], SplitFS [SOSP '19]

Hardware: 8-core Intel Xeon 4215R CPU

1 × 128GB Intel Optane PM

# Evaluation: Concurrent 4 KB Random Read

● MadFS    ▲ ext4-DAX    ■ NOVA    ◆ SplitFS

# Evaluation: Concurrent 4 KB Random Read

● MadFS　　▲ ext4-DAX　　■ NOVA　　◆ SplitFS



100% Read

Best performance: MadFS

Single thread

- 43% faster than ext4-DAX

- 41% faster than NOVA

All FS scale well: no writes

# Evaluation: Concurrent 4 KB Random Overwrite

● MadFS    ▲ ext4-DAX    ■ NOVA    ◆ SplitFS



MadFS doesn't update kernel metadata

Saturates device bandwidth w/ 1 thread

- 26% faster than SplitFS

- 70% faster than ext4-DAX

High throughput w/ more threads

- Lock-free concurrency control

# Evaluation: TPC-C on SQLite

Transaction processing benchmark on relational database

Characteristic: block-aligned writes followed by fsync

# Evaluation: TPC-C on SQLite

*Result*: MadFS outperforms other filesystems

Mix: 26% faster than SplitFS, 58% faster than ext4-DAX

# Evaluation: More in Paper

## Multi-threaded benchmarks

- Contended concurrent write

- Concurrency control comparison

## Metadata operations

- Open latency

- Garbage collection

## Macro-benchmarks

- YCSB on LevelDB

# Conclusion

## Metadata Embedding

- Many data ops coupled with metadata updates ⇒ expensive kernel I/O stack

- Embed metadata into file data for efficient userspace management

## Per-File Virtualization

- Push file functionalities into userspace as much as possible
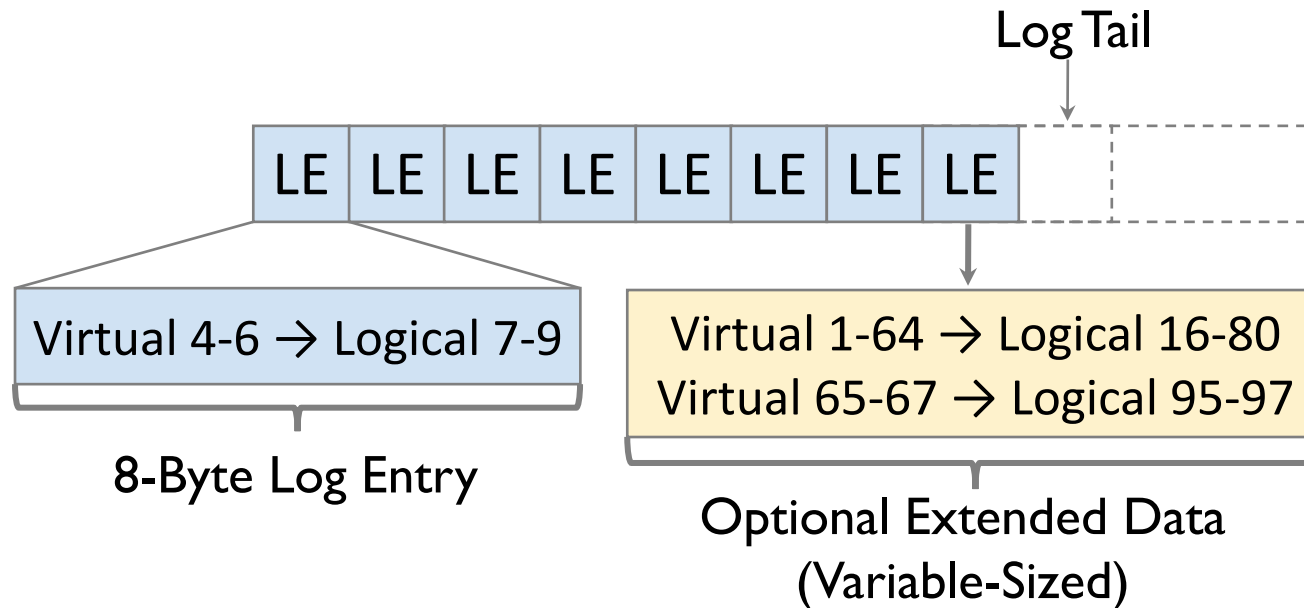
## MadFS: Metadata Embedded Filesystem

- Highly-scalable userspace PM filesystem with strong crash consistency

# Questions

# Backup Slides

# MadFS: Log-Structured Metadata

Fix-sized 8-byte log entries + optional extended data

Log Tail

| LE | LE | LE | LE | LE | LE | LE | LE |

Virtual 4-6 → Logical 7-9

8-Byte Log Entry

Virtual 1-64 → Logical 16-80
Virtual 65-67 → Logical 95-97

Optional Extended Data
(Variable-Sized)

# MadFS: Lock-Free Optimistic Concurrency Control

Use Compare-and-Swap to commit 8-byte log entry

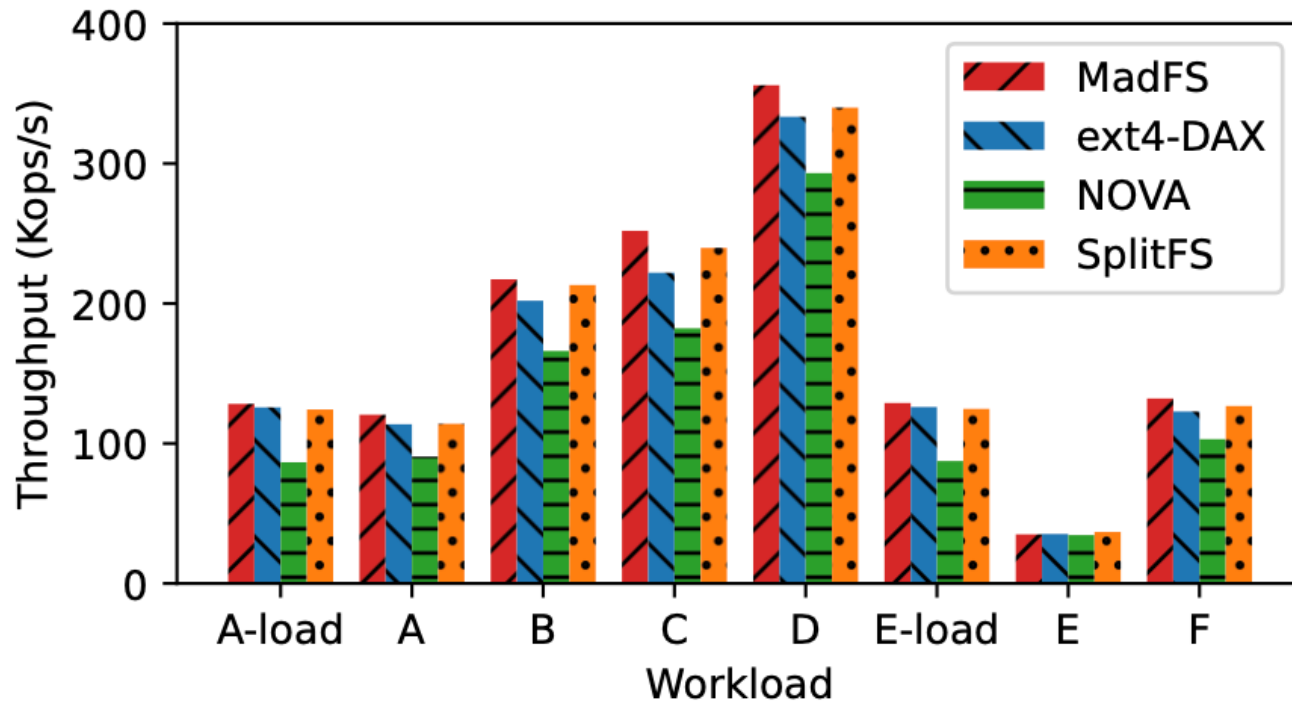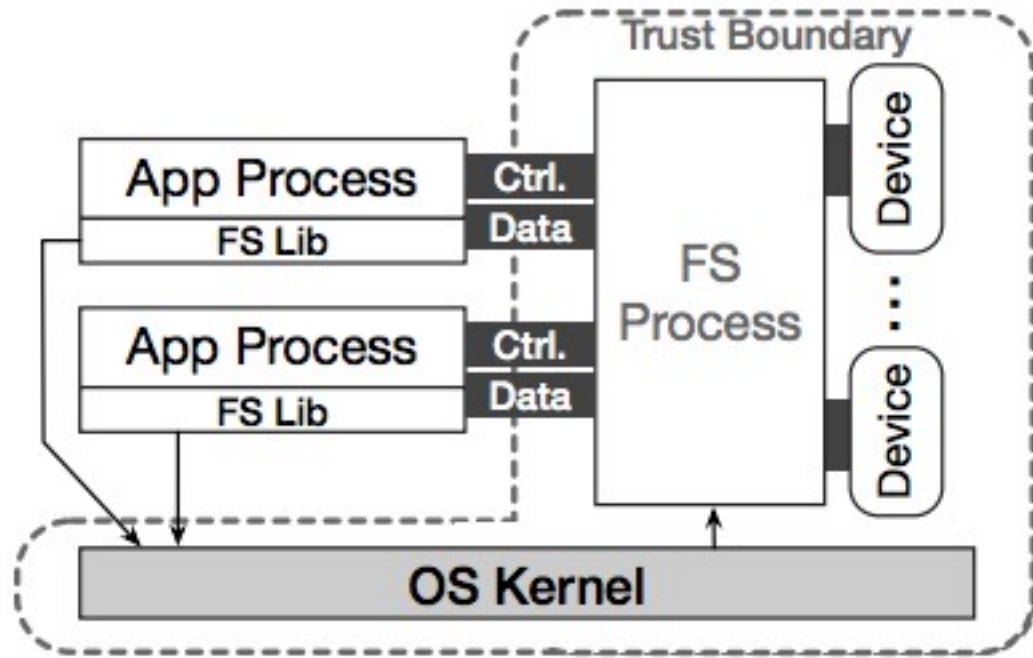| | | | |
|---|---|---|---|
| **1** | **2** | **3** | **4** |
| *Begin* Record Starting Log Tail | *Execute* Copy-on-Write | *Validate* "Compare" if Log Tail still Points to an Empty Slot | *Commit* "Swap" the Log Entry to the Tail |

# Evaluation: YCSB on LevelDB



Workload C

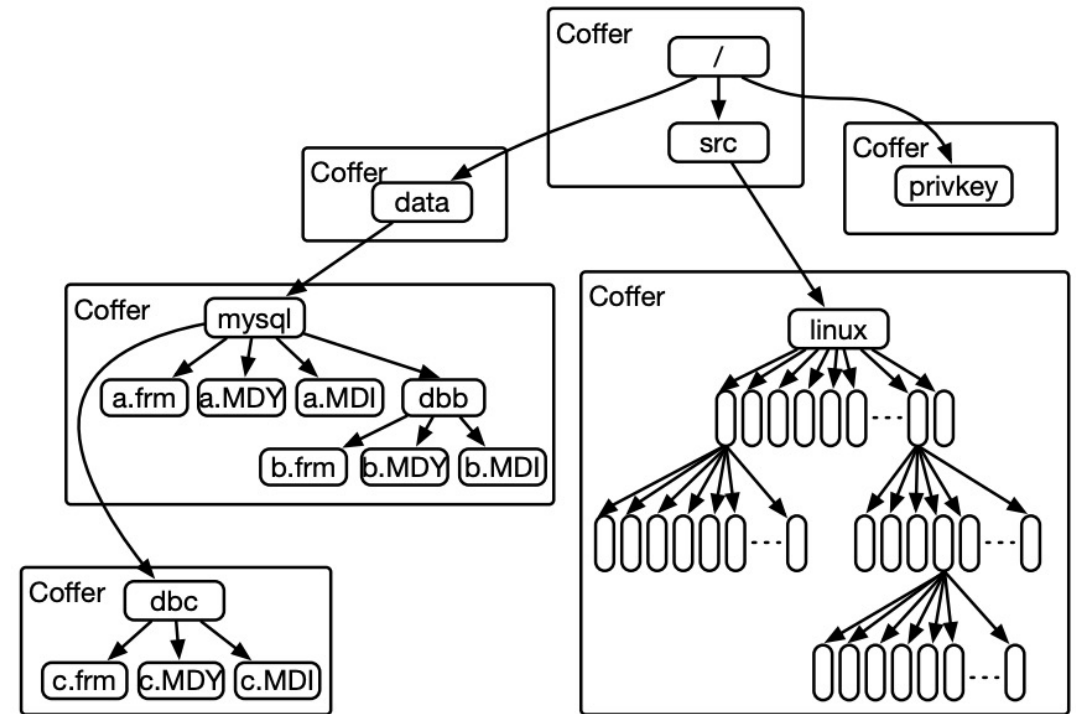- 5% faster than SplitFS

- 12% faster than ext4-DAX

Workload F

- 4% faster than SplitFS

- 7% faster than ext4-DAX

# Related Work



Aerie



ZoFS